

Finite Difference discretization of a simple Convection Problem (ODE): analysis, implementation, verification

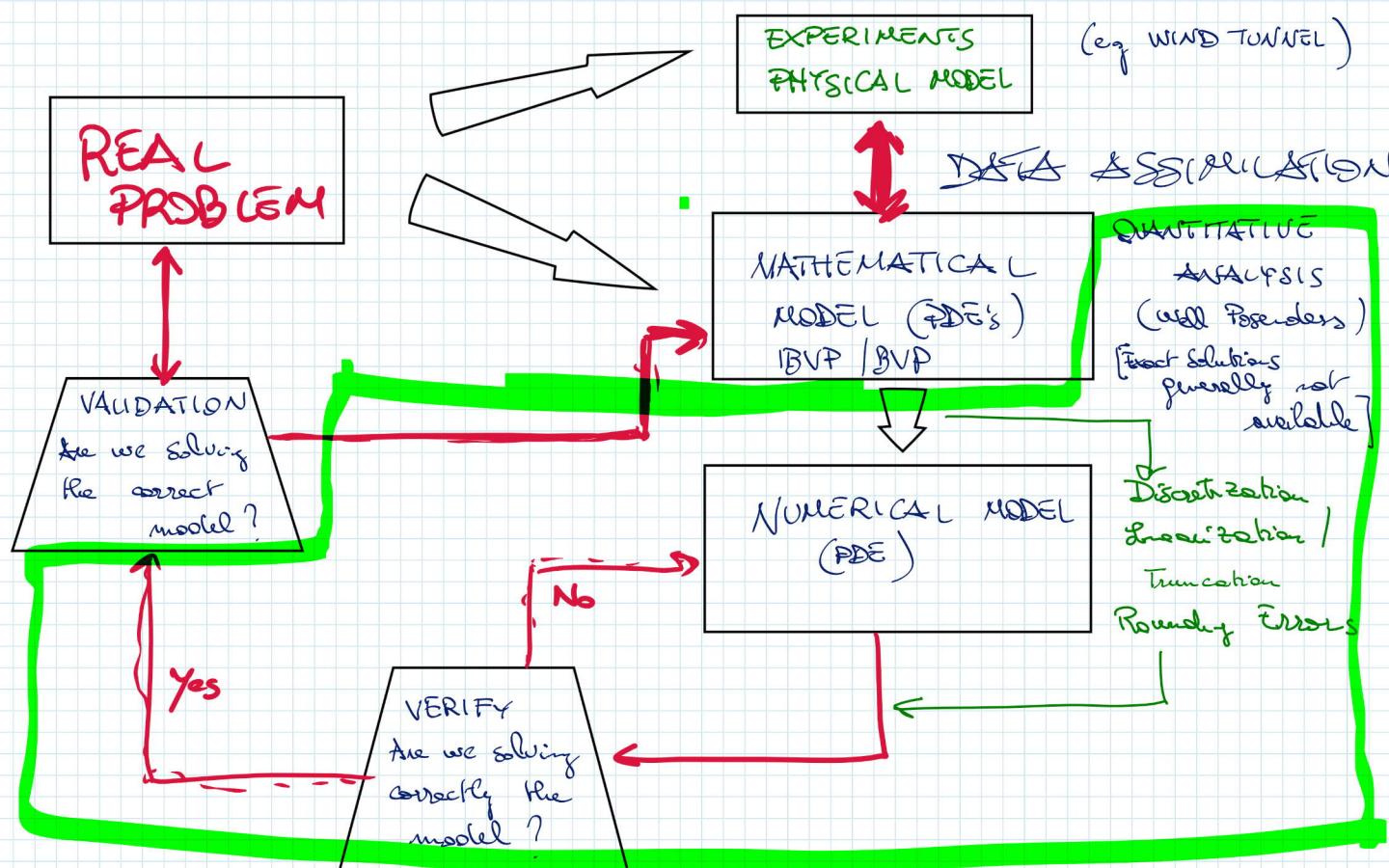
REFERENCE : Langtangen - Linge , Chapter 1 + Appendix B

INTRODUCTION

In this first week, we consider some differential problem numerically approximated by the finite difference method. However, before we consider an example, we introduce some basic concept useful for the entire course. This will help to understand the features we want to investigate in a numerical method for solving a partial differential equation (and beyond).

Let's start with a practical approach leading to the numerical approximation of a Boundary Value Problem (BVP) or an Initial Boundary Value Problem (IBVP).

Modeling a Real problem with Numerical PDE's



In this course we are concerned with the region highlighted in green.

Source of Errors

Discretization

With "discretization" we mean the approximation of the differential components (derivatives) of the problem with algebraic operations. This step is necessary to manage the problem with a finite machine like a computer, yet it introduces some errors.

To give a simple example :

$$u'(x) = \frac{du}{dx}(x) \approx \frac{u(x+h) - u(x-h)}{2h} = (u(x))_h$$

Problem Numerical Problem

The derivative $\left(\lim_{\Delta x \rightarrow 0} \frac{u(x+\Delta x) - u(x)}{\Delta x} \right)$ is approximated by an algebraic operation $(+, -, :, \times)$

Now, reasonable questions are (h = discretization parameter)

⇒ what is the error $|u(\bar{x}) - (u(\bar{x}))_h|$?

More precisely:

- $\lim_{h \rightarrow 0} (u(\bar{x}))_h = u(\bar{x})$?

CONVERGENCE

- if yes: $u(\bar{x}) - (u(\bar{x}))_h = O(h^p)$
what is p

CONVERGENCE ORDER

This is the fundamental feature: WE NEED CONVERGENCE to trust a method

|| This helps ranking different convergent methods: higher p and faster is the method to converge

The analysis of the discretization errors

$|u(\bar{x}) - (u(\bar{x}))_h|$ is one of the specific topics of this course for the different methods proposed.

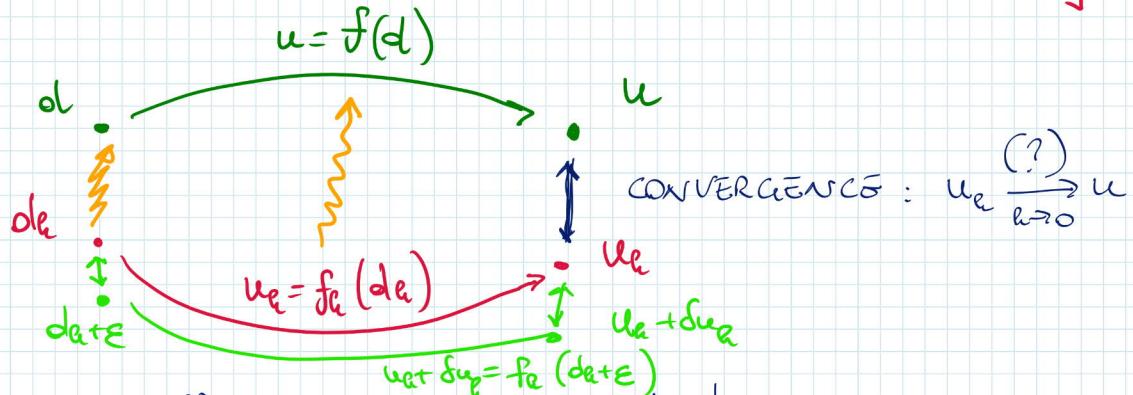
However, let's generalize the ideas:

Mathematical Problem:

$$P(u; d) = 0 \quad \begin{array}{l} u = \text{unknown} \\ d = \text{data, parameters} \end{array}$$

Numerical (Discretized) Problem:

$$P_h(u_h; d_h) = 0 \quad \begin{array}{l} u_h = \text{numerical solution} \\ d_h = \text{approximation of data} \end{array}$$



The convergence generally requires some ingredients

CONSISTENCY : for $h \rightarrow 0$ ($d_h \rightarrow d$)

Then $f_h \rightarrow f$

STABILITY : The numerical problem is "stable" for perturbations:

There exists $C > 0$ s.t.

$$\|\delta u_h\| \leq C\|\epsilon\|$$

with $\|\cdot\|$ proper norms. C in general may depend on h .

(but it is bad if $C \rightarrow +\infty$ when $h \rightarrow 0$)

For linear problems and linear methods (Lax-Richtmyer Theorem)
CONVERGENCE = STABILITY + CONSISTENCY

For nonlinear problems, consistency and stability are only necessary, not sufficient (compactness needed).

In the example above we have that

$$(u'(x))_h = \lim_{h \rightarrow 0} \frac{u(x+h) - u(x-h)}{2h} = u'(x) \text{ by definition (CONVERGENCE)}$$

and with a Taylor expansion we notice that

$$|u'(x) - (u'(x))_h| = O(h^2). \quad (\text{2nd order method}).$$

Notice that in this case consistency and convergence coincide, because the solution of the problem (the derivative) is the problem itself.

However, notice that this analysis does not include the rounding errors.

We will follow-up on this.

For the moment, let's assume that $u(x)$ can be differentiated and that we perturb $u(x)$ with a differentiable function δ s.t. both $\|\delta\|_\infty$ and $\|\delta'\|_\infty$ are bounded \Rightarrow that $\varepsilon = \max(\|\delta\|_\infty, \|\delta'\|_\infty)$.

$$\text{Then } (u+\delta)'_h(x) = \frac{u(x+h) - u(x-h)}{2h} + \frac{\delta(x+h) - \delta(x-h)}{2h} \quad (\text{u}'(x))$$

so we have:

$$|(u+\delta)'_h(x) - (u)'_h(x)| \leq \frac{\varepsilon}{h} \quad \text{for } h > 0 \quad (1)$$

(This estimate doesn't require the differentiability of δ)

If we let $h \rightarrow 0$ we have

$$|(u+\delta)'_h(x) - (u)'_h(x)| \leq |\delta'(x)| \leq \varepsilon$$

(This estimate requires the differentiability of δ)

So, if the perturbation is (reasonably) regular, we have stability (some regularity of the solution)

IN THIS COURSE WE MAINLY INVESTIGATE DISCRETIZATION ERRORS

Linearization / Truncation

As we will see, the final numerical solution may be affected by other errors.

For instance, in most of the cases, we need to solve - after the discretization, a linear system. The numerical solution of this system introduces some errors. In the vast majority of practical problems, the numerical methods for linear systems are ITERATIVE (GMRES, Conjugate Gradient, etc.), so this error is under control with the convergence test (see MATH 517).

Also, when we have a non-linear differential problem, we need to "linearize" it, e.g. with a Newton approach. We have, generally, two options: we DISCRETIZE first and then LINEARIZE (D-L) or the other way around (L-D). The two approaches are not equivalent, D-L is generally preferred. The linearization process introduces other errors, because we need to solve non-linear algebraic equations with some method investigated in the courses of numerical analysis. These errors are not specifically analyzed here, yet they are important in some cases. Since numerical methods for nonlinear problems are iterative, also in this case we can control the error with the stopping criterion.

EXAMPLE :

$$-\frac{\partial^2 u}{\partial x^2} + u \frac{\partial u}{\partial x} + u^3 = f \quad x \in (0,1)$$

$$u(0) = u(1) = 0$$

CONSISTENCY ERROR

With "consistency error", we mean specifically the residual ($=$ quantity $\neq 0$) we obtain when we enforce the exact solution of a problem to the numerical scheme.

For instance, for $u = e^{2x}$, $u'(x) = 2e^{2x}$ and:

$$(u)_h'(x) = \frac{u(x+h) - u(x-h)}{2h} = \frac{e^{2(x+h)} - e^{2(x-h)}}{2h} = e^{2x} + O(h^2)$$

consistency error

A possible DISCRETIZE-then-LINEARIZE approach:

GRID:

$$\Omega = [x_0, x_1, x_2, \dots, x_N] \quad h = \frac{x_N - x_0}{N}$$

Definition:

$$u_i(x_j) = u_j$$

$$\left\{ \begin{array}{l} \frac{u_{j+1} - 2u_j + u_{j-1}}{h^2} + u_j \frac{u_{j+1} - u_{j-1}}{2h} + u_j^3 = f_j \\ u_0 = 0 \quad u_N = 0 \end{array} \right.$$

This is a nonlinear algebraic system that we can write as:
 $\underline{F}(u) = 0$ with
 $u = [u_0, u_1, \dots, u_N]^T$

Newton method:

Denote by J the Jacobian matrix $[J_{ij} = \frac{\partial F_i}{\partial u_j}]$

Then we start with a guess $\underline{u}^{(0)}$ and we start the loop:

$k = 0$

while (Not converge)

Compute $J(\underline{u}^{(k)})$

Solve the linear system:

$$J(\underline{u}^{(k)}) \delta = -\underline{F}(\underline{u}^{(k)})$$

$$\text{Update: } \underline{u}^{(k+1)} = \underline{u}^{(k)} + \delta$$

$k++$

$$\text{Convergence} = \text{Rel-Convergence} (\underline{u}^{(k+1)}, \underline{u}^{(k)}, \underline{F}(\underline{u}^{(k)}))$$

end-while

There is another possible scheme:

$$e \frac{\partial u}{\partial x} = \frac{1}{2} \frac{\partial(u^2)}{\partial x} \rightarrow \frac{1}{2} \frac{u_{j+1}^2 - u_{j-1}^2}{2h}$$

This approach is called CONSERVATIVE

In the Linearize-then-Discretize approach, we need to differentiate the equation with respect to a function (Gâteaux Derivative). Although possible, we do not follow this approach here.

Round-off Errors

All the numbers represented on a computer are subject to numerical approximation. This approximation depends on the hardware (32 bit vs 64 bit for instance) and it is quantified by the MACHINE-E, i.e. the smallest number ϵ such that $1+\epsilon > 1$ for the computer. Generally, $\epsilon \approx 10^{-16}$.

We do not consider the presence of round-off errors when we do our convergence analysis, but we need to keep in mind that round-off errors are everywhere!

EXAMPLE:

$$e^x = \frac{d e^x}{dx} \Big|_{x=\bar{x}} \approx \frac{e^{\bar{x}+h} - e^{\bar{x}}}{h}$$

Theory: for $h \rightarrow 0$, $\lim_{h \rightarrow 0} \frac{e^{\bar{x}+h} - e^{\bar{x}}}{h} = e^{\bar{x}}$

Practice: $f(x) = e^{\bar{x}+h}$ ($=$ floating point approximation of $e^{\bar{x}+h}$) $= e^{\bar{x}+h} (1 + \epsilon_1)$

$$f(x) = e^{\bar{x}} (1 + \epsilon_2)$$

with $|\epsilon_1|, |\epsilon_2| \leq \epsilon$

So:

$$\frac{f(x) - f(x)}{h} = \frac{e^{\bar{x}+h} - e^{\bar{x}}}{h} + \frac{e^{\bar{x}}(e^h - 1)}{h}$$

for $h \rightarrow 0$
then $\rightarrow e^h$

$$\left| \frac{e^{\bar{x}}(e^h - 1)}{h} \right| < \frac{e^{\bar{x}} 2\epsilon}{h} \xrightarrow[h \rightarrow 0]{} 0$$

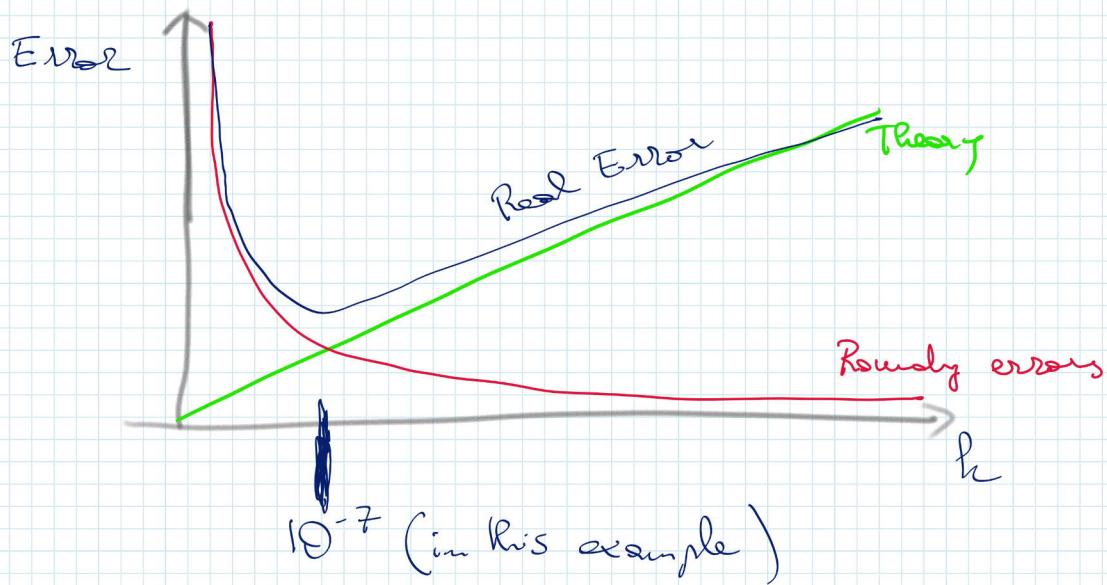
Keep in mind this aspect in your numerical simulations.

TRY: implement this simple example

Python:

```
import numpy as np  
x = 2  
h = 0.1  
np.abs(np.exp(x) - ((np.exp(x+h)-np.exp(x))/h))
```

Test the error for decreasing h : for $h > 10^{-7}$ everything goes as expected (linear reduction of the error), for $h < 10^{-7}$ you see the error going up



Question: So, why we do not consider rounding errors in our analysis?

Answer: in general, practical problems are OK with errors greater than the range of h where rounding errors show up ($h \gg 10^{-7}$).

FINITE DIFFERENCES

The first method we consider is the simplest one: simplicity means that we do not need much background, but also that there are inherent limitations (we will go over when studying other methods). The basic idea is to replace the derivatives with incremental quotients in specific points selected with some criterion. We start with a simple INITIAL VALUE problem, so to avoid some technicalities, yet to get familiar with some concepts and analysis tools.

LESSONS FROM A SECOND ORDER CAUCHY PROBLEM

We start considering the numerical solution of

$$\frac{d^2 u}{dt^2} + \omega^2 u = 0 \quad t > 0$$

$$u(0) = I \quad \frac{du}{dt}(0) = 0$$

Clearly, we have the analytical solution $u(t) = I \cos(\omega t)$. This will help us understand how accurate are the methods we consider.

To perform a numerical approximation, let's "discretize" the differential equation, replacing the derivative with algebraic operations:

$$P(u) = 0 : \frac{d^2 u}{dt^2} + \omega^2 u = 0 ; \quad u(0) = I , \quad u'(0) = 0 .$$

Numerical Approximation



Let's collocate the problem in the points $t_j = jh$ and replace the derivatives with the incremental quotients:

$$(2) \quad \frac{u_{j+1} - 2u_j + u_{j-1}}{h^2} + \omega^2 u_j = 0$$

$u_0 = I$ For $u'(0)$ we need an approximation of the first derivative in 0

Let's start with the option in the book: $u'(0) = \frac{u_1 - u_{-1}}{2h}$ This is clearly a numerical trick, because this point doesn't exist.

However, it can work.

So, we can write sequentially:

$$\begin{aligned} \text{I.C.} & \left\{ \begin{array}{l} u_0 = I \\ \frac{u_1 - u_{-1}}{2h} = u'(0) = 0 \Rightarrow u_{-1} = u_1 \end{array} \right. \\ \text{First Step} & \left\{ \begin{array}{l} \frac{u_1 - 2u_0 + u_{-1}}{h^2} + \omega^2 u_0 = 0 \Rightarrow \frac{2u_1 - 2I}{h^2} + \omega^2 I = 0 \\ u_1 = \left(1 - \frac{\omega^2 h^2}{2} \right) I \end{array} \right. \end{aligned}$$

$$\text{Subsequent steps} \quad \left\{ \begin{array}{l} u_{j+1} = \left(2 - \omega^2 h^2 \right) u_j - u_{j-1}, \quad j \geq 2 \end{array} \right.$$

The method is clear and pretty immediate: how does it perform?

SEE THE PYTHON CODE

ex1.py

If we use the code and we measure the error as:

$$e_{\infty} \stackrel{\text{def}}{=} \max_j |u_{\infty}(t_j) - u_j|$$

we find that $e_{\infty} \sim O(h^2)$.

```
def solver(l, w, dt, T):
```

```
    ...  
    Solve u'' + w**2*u = 0 for t in (0,T], u(0)=l and u'(0)=0,  
    by a central finite difference method with time step dt.  
    ...  
    dt = float(dt)  
    Nt = int(round(T/dt))  
    u = np.zeros(Nt+1)  
    t = np.linspace(0, Nt*dt, Nt+1)  
  
    u[0] = l  
    u[1] = u[0] - 0.5*dt**2*w**2*u[0]  
  
    for n in range(1, Nt):  
        u[n+1] = 2*u[n] - u[n-1] - dt**2*w**2*u[n]  
    return u,
```

i.e.

core of
the method

CONVERGENCE ANALYSIS

Let's try to explain our numerical results.

If we use a Taylor expansion, we have

$$\left\{ \begin{array}{l} u_{j+1} = u_j + u'_j h + u''_j \frac{h^2}{2} + u'''_j \frac{h^3}{3!} + u''''_j \frac{h^4}{4!} + \text{High Order Terms} \\ u_{j-1} = u_j - u'_j h + u''_j \frac{h^2}{2} - u'''_j \frac{h^3}{3!} + u''''_j \frac{h^4}{4!} + \text{H.O.T.} \end{array} \right.$$

If we sum these two, we have:

$$u''_j = \frac{u_{j+1} - 2u_j + u_{j-1}}{h^2} + \frac{u''''_j h^4}{12} + \text{H.O.T.}$$

So, we easily define the LOCAL TRUNCATION ERROR:

$$\text{Exact Problem: } u''_{\infty}(x_j) + w^2 u_{\infty}(x_j) = \frac{u_{\infty,j+1} - 2u_{\infty,j} + u_{\infty,j-1}}{h^2} + w^2 u_{\infty,j} = -\frac{w^2 h^2}{12}$$

$$\text{Numerical Problem: } \frac{u_{j+1} - 2u_j + u_{j-1}}{h^2} + w^2 u_j = 0$$

The Local Truncation Error $\xrightarrow[h \rightarrow 0]{} 0$

This states that the method is consistent with the original problem.

Unfortunately this is necessary but not sufficient for the convergence.

In this simple problem, we can work out the convergence directly.

Set $e_j = u_{\infty,j} - u_j$. Our goal is to prove that the $|e_j| \xrightarrow[h \rightarrow 0]{} 0$.

By subtracting the exact problem and the numerical one and setting the LTE: $\frac{u_{\infty}(t_j)}{12} h^2 = g h^2$ we have:

$$e_{j+1} - (2 - w^2 h^2) e_j + e_{j-1} = g_j h^4$$

$$e_0 = 0$$

$$e_1 = Dh^2$$

$$\Rightarrow \left[\frac{u_1 - u_{-1}}{2h} + Dh^2 = u'(0) \right]$$

By Taylor expansion we can prove this

Let's solve this 2nd order difference equation. Before, a REMARK

REMARK : On the Taylor Expansion (TE)

When we do TE, we postulate "regularity", i.e. that all the derivatives we need exist. In this case, this is true because $u_{\infty} = T \cos(wt)$ is infinitely regular. However, when we don't know the solution, LACK OF REGULARITY may undermine our analysis. This may be the root of discrepancy between the expectations (theory) and the results (practice).

If we take $C_j \approx C$, we have an error equation in the form: (C can be $C = \max_i C_i$)

$$e_{j+1} - (2 - \omega^2 h^2) e_j + e_{j-1} = Ch^4$$

(*) Difference equation of order 2; please check how to solve it.

A particular solution to this difference equation is for $e = \text{constant}$:

$$e = \frac{Ch^4}{1 - 2 + \omega^2 h^2 + 1} \approx \frac{C}{\omega^2} h^2$$

Then, we add errors from the homogeneous part: to do this, we need to solve:

$$\gamma^2 - (2 - \omega^2 h^2) \gamma + 1 = 0 \Rightarrow \gamma_{1,2} = \frac{2 - \omega^2 h^2 \pm \sqrt{A + \omega^4 h^4 - 4\omega^2 h^2}}{2} = \frac{2 - \omega^2 h^2 \pm \omega h \sqrt{\omega^2 h^2 - 4}}{2}$$

Then, the solution will be:

$$(a) \text{ for } \gamma_1 \neq \gamma_2: e_j = A \gamma_1^j + B \gamma_2^j + \frac{C}{\omega^2} h^2 \quad \text{where } A \text{ and } B \text{ depend on the initial conditions.}$$

In our case, $e_0 = 0 \quad e_1 = Dh^2$ (D is another constant).

In fact $\begin{cases} u_{\alpha}(0) = 0 = \frac{u_1 - u_{-1}}{2h} + Dh^2 \\ u'(0) = \frac{u_1 - u_{-1}}{2h} \end{cases}$

check with Taylor expansion

$$(b) \text{ for } \gamma_1 = \gamma_2 = \gamma \quad e_j = (A + tB) \gamma^j + \frac{C}{\omega^2} h^2$$

Let's start considering the case:

$$(1) \quad \gamma_1 \neq \gamma_2 \text{ real: this occurs for } \omega^2 h^2 - 4 > 0 \Rightarrow h > \frac{2}{\omega} \quad (\text{a time step larger than the inverse of half of the frequency})$$

$$\gamma_{1,2} = \frac{2 - \omega^2 h^2 \pm \omega h \sqrt{\omega^2 h^2 - 4}}{2} \Rightarrow \text{Notice that } \gamma_1 \gamma_2 = 1, \text{ so } \gamma_1 = \frac{1}{\gamma_2}. \quad \text{If } |\gamma_2| < 1, \text{ then } |\gamma_1| > 1 \quad (\text{or the other way around})$$

so either γ_1 or γ_2 is $|\gamma_j| > 1$

This means $|e_j| \xrightarrow{j \rightarrow \infty} +\infty$

We do not have convergence because of the error accumulation

$$(2) \quad \gamma_1 = \gamma_2: \quad h = \frac{2}{\omega} \Rightarrow \gamma = \frac{2 - \omega^2 h^2}{2} = -1$$

$$e_j = (A + jB)(-1)^j + \frac{Ch^2}{\omega^2} \Rightarrow e_0 = A + C \frac{1}{\omega^4} = 0 \rightarrow A = -\frac{C}{\omega^4}$$

$$\Rightarrow e_1 = \frac{Ch^2}{\omega^2} - B = Dh^2$$

if $B = 0$, we expect an oscillatory behavior because of the $(-1)^j$ term.

if $B \neq 0$, we may expect an explosion

So, in general $h = \frac{2}{\omega}$ can be unstable

$$(3) \quad \gamma_1 = \gamma_2 \text{ complex conjugate: } h < \frac{2}{\omega}$$

$$\gamma_{1,2} = \frac{2 - \omega^2 h^2 \pm i\omega h \sqrt{4 - h^2 \omega^2}}{2}$$

$$\text{Notice that: } |\gamma_j| = \frac{1}{2} \sqrt{(4 + \omega^4 h^4 - 4\omega^2 h^2 + 4\omega^2 h^2 - \omega^4 h^4)} = 1$$

So we can conclude that:

$$|e_j| \leq |A||\gamma_1|^j + |B||\gamma_2|^j + \left| \frac{Ch^2}{\omega^2} \right| \leq A + B + \frac{Ch^2}{\omega^2}$$

$$0 = e_0 = A + B + \frac{Ch^2}{\omega^2}$$

$$Dh^2 = e_1 = A\gamma_1 + B\gamma_2 + \frac{Ch^2}{\omega^2}$$

$$\begin{aligned} A + B &= -\frac{Ch^2}{\omega^2} \\ A(\gamma_1 - 1) + B(\gamma_2 - 1) &= Dh^2 \\ A(\gamma_1 - 1) &= \left(D + \frac{(\gamma_2 - 1)C}{\omega^2} \right) h^2 \\ B(\gamma_2 - 1) &= \left(D + \frac{(\gamma_1 - 1)C}{\omega^2} \right) h^2 \end{aligned} \Rightarrow |A|, |B| \sim O(h^2)$$

We conclude that : for $h < \frac{2}{\omega}$ the method is consistent and stable, so it converges and, in particular, with order 2.

REMARK

Notice that the convergence rate is dictated by the consistency error for the equation, but also for the initial state.

EXERCISE 1:

Approximate the initial derivative with no ghost point : $u'(0) \approx \frac{u_1 - u_0}{h} \Rightarrow u_1 = u_0$

Test and prove that the convergence rate becomes 1 in this case.

EXERCISE 2 :

If we want a second order formula, with no ghost point, we can use the formula :

$$u'(0) \approx \frac{3u_0 - 4u_1 + u_2}{2h}$$

(1) Prove that this is a 2nd order approximation of $u'(0)$

(2) Test the method with this formula

(3) Prove that the method converges with order 2.

ANOTHER (MORE SPECIFIC) CONVERGENCE ANALYSIS

The solution to this problem is PERIODIC. When we think to the Fourier decomposition, we know that we can do a spectral or frequency analysis.

$$f \text{ periodic or on a bounded interval} = \sum_{j=0}^{\infty} \hat{f}_j e^{i j t}$$

In our particular case, the solution is a co-sinusoidal function $u_\infty = I \cos(\omega t)$.

In this case, we may be interested specifically to the errors on the amplitude I and the frequency ω .

In other terms assuming that the numerical solution reads like : $u_j = I_j \cos(\omega_j t_j)$

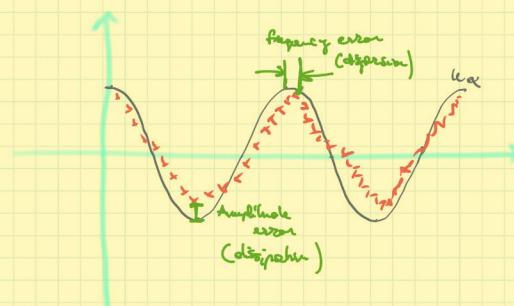
where we have a numerical amplitude I_j and a numerical frequency ω_j , what are the :

(1) AMPLITUDE ERROR ($|I - I_j|$)

(2) FREQUENCY ERROR $|\omega - \omega_j|$ or $\frac{\omega}{\omega_j}$?

Amplitude errors go under the name of DISSIPATION ERRORS

Frequency errors go under the name of DISPERSION ERRORS



We can categorize the error according to these two components (Von Neumann analysis).

REMARK : The following analysis is equivalent but different of the book. It is more general

Let's compute directly the numerical solution and see if we can identify AMPLITUDE I_j and FREQUENCY ω_j .

Based on this, we can compute specifically the different components of the errors.

The numerical solution solves :

$$u_{j+1} - (2 - \omega^2 h^2) u_j + u_{j-1} = 0$$

(I denote ω_j as $\tilde{\omega}$ for easier of notation)

Based on what we learned with the error equation, for $h < \frac{2}{\omega}$ we have :

$$u_j = A \tilde{s}_1^j + B \tilde{s}_2^j$$

$$\tilde{s}_{1,2} = \frac{2 - \omega^2 h^2}{2} \pm i \frac{\sqrt{4 - \omega^2 h^2}}{2}$$

$$\text{Notice that: } |\tilde{s}| = \sqrt{4 - \omega^2 h^2} = 1$$

$$\text{Using the polar form: } \tilde{s}_1 = |\tilde{s}| e^{i \tilde{\omega} h} \quad \tilde{s}_2 = |\tilde{s}| e^{-i \tilde{\omega} h} \quad \text{where} \quad \tilde{\omega} = \frac{1}{h} \arctan \frac{\sqrt{4 - \omega^2 h^2}}{2 - \omega^2 h^2}$$

For simplicity, set $y = \frac{\omega h}{2 - \omega^2 h^2}$

Now, let us compute A and B: in our scheme:

$$u_0 = I = A + B$$

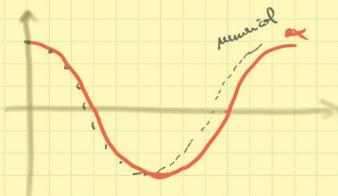
$$u_i - u_{i-1} = 0 \Rightarrow Ae^{i\omega h} + Be^{-i\omega h} = Ae^{i\omega h} + Be^{i\omega h} \Rightarrow A(e^{i\omega h} - e^{-i\omega h}) = B(e^{i\omega h} - e^{-i\omega h}) \Rightarrow A = B$$

$$u_j = \frac{I}{2} (e^{i\tilde{\omega} h j} + e^{-i\tilde{\omega} h j}) = I \cos(\tilde{\omega} h j)$$

The exact solution is $u_\alpha(t_j) = I \cos(\omega h j)$

From here it is apparent that (in exact arithmetic) $I_\alpha = I$ and $\omega_\alpha = \tilde{\omega} = \frac{1}{h} \arctan(y)$

We do not expect dissipation error, we also expect dispersion errors.
This is consistent with the numerical results.



Let's investigate $\frac{1}{h} \arctan(y)$

(1°) \arctan is an odd function. For $h \rightarrow 0$, $y \rightarrow 0$ and

$$\arctan(y) \approx y - \frac{1}{3} y^3 + \frac{1}{5} y^5$$

(2°) $y = \frac{wh \sqrt{4 - \omega^2 h^2}}{2 - \omega h^2} = \frac{\sqrt{4 - \omega^2}}{2 - \omega h^2}$ \Rightarrow This is odd in x too, so for $x \rightarrow 0$ (even derivatives are 0)

$$y(x) \approx x + \frac{3}{8} x^3 + c x^5$$

Putting together (1°) and (2°) we have: $\arctan(y(x)) \approx y - \frac{1}{3} y^3 \approx x + \frac{3}{8} x^3 - \frac{1}{3} x^3 = wh \left(1 + \frac{1}{24} \omega^2 h^2\right)$

$$\Rightarrow \tilde{\omega} = \omega \left(1 + \frac{1}{24} \omega^2 h^2\right) (+ O(h^4))$$

This is the same conclusion of the book:

$$\text{The dispersion error } \frac{\tilde{\omega}}{\omega} = \left(1 + \frac{\omega^2 h^2}{24}\right) \text{ so (1) } \tilde{\omega} > \omega \text{ (2) } \tilde{\omega} \rightarrow \omega \text{ with } h^2.$$

This analysis will subject of Homework!

Postprocessing

Suppose that we are interested in $u'(x)$. The exact derivative is $-\omega I \sin(\omega t)$.

Once we have u_j we can use a formula like:

$$u'_j = \frac{u_{j+1} - u_{j-1}}{2h} \quad \parallel \text{We know this is a 2nd order formula for the exact values of } u_j.$$

$$u'_\alpha(t_j) = \frac{u_{\alpha,j+1} - u_{\alpha,j-1}}{2h} + O(h^2)$$

$$u'_j = \frac{u_{j+1} - u_{j-1}}{2h} = \frac{u_{\alpha,j+1} - e_{j+1} - u_{\alpha,j-1} + e_{j-1}}{2h} = \frac{u_{\alpha,j+1} - u_{\alpha,j-1}}{2h} - \frac{e_{j+1} - e_{j-1}}{2h} = u'_\alpha(t_j) + O(h^2) + \frac{e_{j+1} - e_{j-1}}{2h}$$

We need to understand how $\frac{e_{j+1} - e_{j-1}}{2h}$ scales with h .

$$e_{j+1} = A\delta_1^j + B\delta_2^j + Ch^2 \quad \text{with} \quad |\delta_{1,2}|=1 \quad \text{and} \quad A, B \sim O(h^2)$$

$$e_{j+1} - e_j = A(\delta_1^{j-1})\delta_1^j + B(\delta_2^{j-1})\delta_2^j = O(h^3) \quad \left. \right\}$$

By direct computation $O(h)$

$$\frac{e_{j+1} - e_j}{2h} \sim O(h^2)$$

We can verify this numerically (hands-on session)

REMARK: Pay Attention to the different ways of implementing the computation of the derivative.

A Different approach for the Numerical Discretization

From the analytical point of view, we can reduce the problem to a first order set of equations with the introduction of the auxiliary variable v :

$$\begin{aligned} u'' + \omega^2 u &= 0 \\ u(0) &= I \\ u'(0) &= 0 \end{aligned} \quad \iff \quad \begin{cases} u' = v \\ v' = -\omega^2 u \\ u(0) = I \\ v(0) = 0 \end{cases} \quad \left. \right\}$$

We can design now different methods with different approximations for u and v with 1st order methods.

When we get to first order problems, we have a lot of options: see MATH516

Quick Recap:

$$\frac{dy}{dt} = f(y) \quad , \quad y(0) = y_0$$

→ One step linear methods:

$$y^{n+1} = y^n + h \left(\theta f(t^n, y^n) + (1-\theta) f(t^{n+1}, y^{n+1}) \right)$$

for $\theta=0$: Explicit/Forward Euler, conditionally stable, First Order

for $\theta=1$: Implicit/Backward Euler, unconditionally stable, First Order

for $\theta=\frac{1}{2}$: Crank-Nicolson, unconditionally stable, Second Order

For $\theta \neq 0$ we need to solve at each step a system (nonlinear if f is nonlinear, linear otherwise)

→ Linear Multi-Step Methods (Adams: Adams-Basforth / Adams-Moulton, BDF)

→ Runge-Kutta : e.g., Heun:

$$y^{n+1} = y^n + \frac{h}{2} \left(f(t^n, y^n) + f(t^{n+1}, y^{n+1}) \right)$$

In our case:

$$f = A\bar{y} = \begin{bmatrix} 0 & 1 \\ -\omega^2 & 0 \end{bmatrix} \bar{y} \quad \text{with} \quad \bar{y} = \begin{bmatrix} u \\ v \end{bmatrix}$$

Notice that A has two distinct eigenvalues: $\lambda = \pm i\omega$ expected for an oscillatory solution

For explicit Euler, we have:

$$y^{n+1} = y^n + h \Delta y^n = (I + hA)y^n = \begin{bmatrix} 1 & h \\ -\omega^2 & 1 \end{bmatrix} y^n.$$

$$\text{Eigenvalues: } (\lambda - 1)^2 + \omega^2 h^2 = 0 \quad \lambda^2 - 2\lambda + 1 + \omega^2 h^2 = 0 \quad \lambda_{1,2} = 1 \pm i\omega h \quad |\lambda| = 1 + \omega^2 h^2$$

Notice also that:

$$y^{n+1} = y^n + h \Delta y^n \Rightarrow \begin{cases} u^{n+1} = u^n + h v^n \\ v^{n+1} = -\omega^2 h + u^n + v^n = -\omega^2 h + (u^n - u^n) / h \end{cases} \Rightarrow u^{n+2} = u^{n+1} + h v^{n+1} = u^{n+1} - \omega^2 h u^n + u^{n+1} - u^n = 2u^{n+1} - (1 + \omega^2 h^2) u^n$$

In summary, the method reads (adjusting the index)

$$u^{n+1} = 2u^n - (1 + \omega^2 h^2) u^{n-1}$$

$$(u^n = (2 - \omega^2 h^2) u^n - u^{n-1})$$

The associated difference equation reads

$$g^2 - 2g + (1 + \omega^2 h^2) = 0 \quad g_{1,2} = 1 \pm i\omega h^2 \quad |g_{1,2}| = 1 + \omega^2 h^2$$

This method for this problem is unable of containing the error for the amplitude $|g_j| \rightarrow \infty \quad \forall h > 0$.

A similar analysis can be done for the other methods.

It is interesting to note that we can see:

$$\text{Explicit Euler: } u^{n+1} = 2u^n - u^{n-1} (1 + \omega^2 h^2)$$

$$\text{Implicit Euler: } u^{n+1} = -\omega^2 h u^{n+1} + 2u^n - u^{n-1}$$

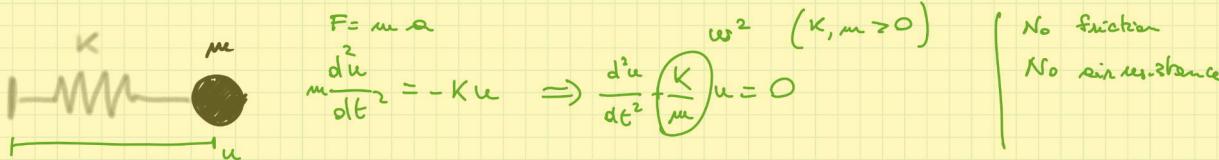
$$\text{Crank-Nicolson: } u^{n+1} = 2u^n - u^{n-1} - \frac{\omega^2 h^2}{4} (u^{n+1} + 2u^n + u^{n-1})$$

$$\text{Our previous method: } u^{n+1} = 2u^n - u^{n-1} - \omega^2 h u^n$$

$$\text{Notice that: } \frac{u^{n+1} + 2u^n + u^{n-1}}{4} \approx u^n + O(h^2)$$

Another Way for Measuring the Accuracy (Based on Physics).

The problem we are considering has a clear physical interpretation:



The energy of the system is the kinetic energy of the mass $(\frac{1}{2}mu^2)$ + the elastic energy of the spring $\frac{1}{2}Ku^2$. For our "normalized" equation we can compute the energy with the virtual work principle.

$$E = \frac{1}{2}(u')^2 + \frac{1}{2}\omega^2 u^2$$

$$\frac{dE}{dt} = \frac{dE}{du} \frac{du}{dt} = u' u'' + \omega^2 u u' = u'(u'' + \omega^2 u) = 0$$

Not surprisingly, the energy of the system is constant:
 $E(t) = E(0) = \frac{1}{2}0 + \frac{1}{2}\omega^2 I^2 = \frac{1}{2}\omega^2 I^2$

A way to assess the performance/accuracy of a method is to compute the numerical energy E_j :

$$E_j = \frac{1}{2} \left(\frac{u_{j+1} - u_{j-1}}{2h} \right)^2 + \frac{1}{2} \omega^2 u_j^2$$

If $|E_j - \frac{1}{2} \omega^2 I^2|$ is small, we can conclude that the method is accurate.

EXERCISE: Test this on the methods considered here

Summary: What we Learned from this example ?

- (1) There are many options in the discretization of an equation: use your fantasy
- (2) Consistency is not enough, we need (at least) STABILITY.
- (3) Direct error analyses are possible and necessary for a full educated interaction with your computer
- (4) For propagative dynamics, the Fourier analysis helps understanding the error
- (5) Analytical solutions are crucial for assessing the performances of a method

FINITE DIFFERENCES FOR A FIRST BOUNDARY VALUE PROBLEM (Elliptic problem)

Ref.: Quarteroni / Sacco / Saleri
NUMERICAL MATHEMATICS

We consider here how our numerical approximation may work when we consider a boundary value problem, a classical one:

$$-K u'' = f(x) \quad x \in (0,1)$$

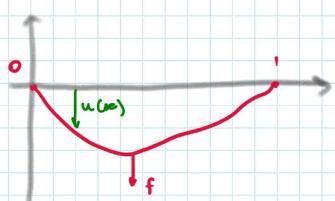
This is the problem of an ^{elastic} string displacement under the force f when small deformations are assumed.

The analytical Problem

Before we give details on the numerical approximation, let's see some insights on the nature of this problem.

Assume that we want to minimize the following functional

$$\bar{J}(u) = \frac{1}{2} \int_0^1 K(u')^2 - \int_0^1 f u \quad \text{with } u(0)=u(1)=0$$



This represents the total energy of a system where the first term is the potential elastic energy and the second term is the energy ^{to win} the external force f to generate the displacement u .
(Let's assume $K > 0$ constant).

To compute the minimum of \bar{J} we compute

$$\frac{\bar{J}(u+\epsilon v) - \bar{J}(u)}{\epsilon}$$

where v is a function with the same regularity of u and $v(0)=v(1)=0$.

$$\frac{\frac{1}{2} K \left(\int_0^1 u'^2 + 2 \epsilon u v + \epsilon^2 v^2 \right) - \int_0^1 u - \epsilon \int_0^1 f v - \frac{1}{2} \epsilon K \int_0^1 u'^2 + \int_0^1 f u}{\epsilon}$$

$$= \int_0^1 u' v' - \int_0^1 f v + \frac{\epsilon}{2} K \int_0^1 (v')^2$$

$$\text{Now, compute } \frac{\partial \bar{J}}{\partial u} \stackrel{\text{def}}{=} \lim_{\epsilon \rightarrow 0} \frac{\bar{J}(u+\epsilon v) - \bar{J}(u)}{\epsilon} = K \int_0^1 u' v' - \int_0^1 f v$$

Virtual Work principle:

The solution of energy minimization solves the problem (with $u(0)=u(1)=0$)

$$K \int_0^1 u' v' = \int_0^1 f v \quad \forall v \in L^2 \quad v(0)=v(1)=0$$

Now, we have,

$$K \int_0^1 u' v' = K \int_0^1 u v' - K \int_0^1 u'' v \stackrel{\text{B.P.}}{=} 0 \quad v(0)=v(1)=0$$

if u is regular enough (u'' exists) then the problem reads

$$\int_0^1 (-K u'' - f) v = 0$$

v is arbitrary

$$-K u'' = f \quad \text{for } x \in (0,1).$$

In other terms, the equation originates from a MINIMIZATION process or, more in general, from a VARIATIONAL procedure (i.e. finding the solution as $\bar{J}(u) \leq \bar{J}(u+\text{variation})$).

This problem will be proved to be well posed later. For now, trust me.

The analytical solution can be found by double integration:

$$-K u'' = f \Rightarrow u''(z) = -\frac{1}{K} \int_0^z f + C \Rightarrow u(x) = -\frac{1}{K} \int_0^x \int_0^z f(z) dz + Cx + D$$

\tilde{z}, z = dummy integration variables,
 C, D = integration constants

With the boundary conditions: $u(0)=0 \Rightarrow D=0$

$$u(1) = -\frac{1}{K} \int_0^1 \int_0^z f(z) dz + C = 0 \Rightarrow u(x) = \frac{1}{K} \left(\int_0^x F(\tilde{z}) d\tilde{z} - \int_0^x F(z) dz \right)$$

$F(z) = \int_0^z f(\tilde{z}) d\tilde{z}$
(anti-derivative of f)

Note that $\int_0^x F(z) dz = \int_0^x F(z) \cdot 1 dz \stackrel{\text{B.P.}}{=} \left[z F(z) \right]_0^x - \int_0^x z f(z) dz = \int_0^x f(z) dz - \int_0^x z f(z) dz$ $u(x) = \frac{x}{K} \int_0^x (1-z) f(z) dz - \int_0^x z f(z) dz$

$\int_0^x F(z) dz \stackrel{\text{B.P.}}{=} \left[z F(z) \right]_0^x - \int_0^x z f(z) dz = x \int_0^x f(z) dz - \int_0^x z f(z) dz$ (we can change name to dummy variable) $= \int_0^x \frac{x-z}{K} f(z) dz + \int_0^x \frac{z-x}{K} f(z) dz = \int_0^x g(x,z) f(z) dz$

where $G(x, z) = \begin{cases} \frac{x(z-x)}{K} & \text{for } 0 \leq z \leq 1 \\ \frac{z(1-x)}{K} & \text{for } 0 \leq x \leq z \end{cases}$ is called Green's function.

If $f \in C^0([0, 1])$, then we have an explicit formula for the solution. Also, since $u'' \propto f \Rightarrow u'' \in C^2([0, 1])$

Remarks:

(1) By a direct computation, one can see that

$$\int_0^1 G(x, z) dz = \frac{1}{2}x(1-x) \Rightarrow \max_{0 \leq x \leq 1} \int_0^1 G(x, z) dz = \frac{1}{8} \Rightarrow \|u\|_\infty \leq \frac{1}{8} \|f\|_\infty$$

(continuous dependence
on the data).

(2) If f is non-negative, then u is non-negative too. This means that the minimum is on the boundary, as easily inferred by the fact that $-u'' = f \Rightarrow u'' = -f \leq 0$

Similarly, if f is non-positive, the maximum is on the boundary (If $f=0$, both max and min are on the boundary). $\left(u = u^{(0)}\right)$

The maximum principle holds also in multiple dimensions

Finite Difference Numerical Approximation

As in the previous chapter/week, we perform the following steps:

(1) Split $[0, 1]$ in subintervals with length h , so that $h = \frac{1}{n}$ $n = \# \text{ of intervals}$ and the vertices are $x_j = jh \quad j = 0, 1, \dots, n$

(2) Collocate the problem at x_j : $-u''(x_j) = f(x_j)$ (Set $K=1$ for simplicity)

(3) Discretize the derivative, e.g.: $-\frac{u_{j+1} - 2u_j + u_{j-1}}{h^2} = f_j \quad u_j \approx u(x_j), \quad f_j = f(x_j)$
with $u_0 = u_n = 0$

In this case, we cannot proceed in a specific direction finding u_{j+1} from u_j and u_{j-1} . But if we write all the equations together, we have a linear system with the matrix:

$$\frac{1}{h^2} \begin{bmatrix} 2 & -1 & 0 & \dots & 0 \\ -1 & 2 & -1 & \dots & 0 \\ 0 & -1 & 2 & -1 & \dots & 0 \\ 0 & 0 & -1 & 2 & -1 & \dots & 0 \\ \vdots & & & & & & \vdots \\ \dots & 0 & 0 & -1 & 2 & \dots & 0 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ \vdots \\ u_{n-1} \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ \vdots \\ f_{n-1} \end{bmatrix} \quad \text{i.e. } A\bar{u} = \bar{f} \quad (\text{with } u_0 = u_n = 0)$$

With the notation introduced earlier:

$u = P(a)$ is the original problem $-u'' = f$, $u(0) = u(1) = 0$

$u_p = P_p(a)$ is the numerical problem $A\bar{u} = \bar{f}$, $u_0 = u_n = 0$

We found that the original problem is well-posed.

What about the numerical one?

The numerical problem is a linear system

$$\underline{A}\underline{u} = \underline{f}$$

$$\underline{A} = \frac{1}{h^2} \text{tridiag}(-1, 2, -1)$$

The question is: is this matrix invertible? We can answer in different ways. Let's use a straightforward one.

Proposition: \underline{A} is symmetric positive definite, so it is invertible

Proof: The symmetry is trivial.

To prove that the matrix is p.d. we can use the following method: a p.d. matrix is such that $\underline{x}^T \underline{A} \underline{x} > 0 \quad \forall \underline{x} \neq \underline{0}$

In our case: $\underline{x}^T \underline{A} \underline{x} = \frac{1}{h^2} \left(2x_1^2 - 2x_1 x_2 + x_2^2 - 2x_2 x_3 + x_3^2 + \dots - 2x_{n-1} x_n + x_n^2 \right) = \frac{1}{h^2} \left(x_1^2 + (x_1 - x_2)^2 + (x_2 - x_3)^2 + \dots + (x_{n-1} - x_n)^2 + x_n^2 \right) > 0 \quad \forall \underline{x} \neq \underline{0}$

(Notice that we can set $x_1 = x_2 = \dots = x_n$ so that we get $\underline{x}^T \underline{A} \underline{x} = \frac{1}{h^2} (x_1^2 + x_n^2) = \frac{2}{h^2} x_1^2$ if we take $x_1 = 0$ then $x_n = 0$ and $\underline{x} = \underline{0}$)

A s.p.d. matrix is invertible (all the eigenvalues are > 0) so we have:

$$\underline{u}_h = \underline{A}^{-1} \underline{f} \quad \text{well posed.}$$

Convergence Analysis

Based on what we did in the previous chapter, we have:

$$-\frac{u_{j+1} - 2u_j + u_{j-1}}{h^2} = -u''(x_j) + ch^2$$

so we can write an error system:

$$\begin{aligned} -u''_h &= -\frac{u_{j+1} - 2u_j + u_{j-1}}{h^2} + ch^2 = f(x_j) \\ -\frac{u_{j+1} - 2u_j + u_{j-1}}{h^2} &= f(x_j) \end{aligned}$$

This states the consistency

$$\Rightarrow \underline{\epsilon} = \underline{u}_h - \underline{u} \quad : \quad \underline{A} \underline{\epsilon} = -\underline{c} h^2$$
$$\underline{c} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix}$$

Vector of the exact values in the nodes:

The question of the convergence refers to the stability. With linear systems, the stability is related to roundings

$$\left. \begin{array}{l} \underline{A} \underline{x}_1 = \underline{b} \\ \underline{A} (\underline{x} + \delta \underline{x}) = \underline{b} + \delta \underline{b} \end{array} \right\} \underline{A} \delta \underline{x} = \delta \underline{b} \Rightarrow \delta \underline{x} = \underline{A}^{-1} \delta \underline{b} \Rightarrow \|\delta \underline{x}\| \leq \|\underline{A}^{-1}\| \|\delta \underline{b}\|$$

We have $\|\underline{A}^{-1}\|$, so we have stability. However, the $\|\underline{A}^{-1}\|$ may depend on h , so the convergence rate is not necessarily $O(h^2)$. Also, consider that if the matrix is close to be singular, the stability may be true but useless. In this case, we are in a good shape.

Proposition : The smallest eigenvalue of \mathbf{A} is independent of h

Proof

It's technical. Let's consider the matrix $T = \text{tridiag}((-1, 2, -1))$. The eigenvalues of \mathbf{A} scales $\frac{1}{h^2}$ * eigenvalues of T . We can compute directly the eigenvalues noting that $T\mathbf{x} = \lambda \mathbf{x}$ needs

$$x_{j+1} - (2-\lambda)x_j + x_{j-1} = 0 \quad x_0 = x_m = 0$$

Solving this, we have:

$$\lambda^2 - (2-\lambda)\lambda + 1 = 0 \quad \lambda_{1,2} = \frac{(2-\lambda) \pm \sqrt{4+\lambda^2-4\lambda-4}}{2} = \frac{(2-\lambda) \pm \sqrt{\lambda^2-4\lambda+8}}{2}$$

If $\lambda > 4$:

$$x_j = A s_1^j + B s_2^j \quad A+B=0 \\ A(s_1^n - s_2^n) = 0 \Rightarrow A=0 \quad (\text{for a generic } n \text{ with } s_1, s_2 \text{ and})$$

So, NO EIGENVALUES > 4 .

If $\lambda < 4$: $\lambda_{1,2} = \frac{2-\lambda \pm \sqrt{\lambda^2-4\lambda+8}}{2}$ Notice that $|s|=1$, so we can write $s_{1,2} = \cos(\theta) \pm i \sin(\theta)$

To have eigenvalues, we need that $\lambda \neq 0$. When we prescribe the boundary conditions, we have:

$$\left. \begin{array}{l} A+B=0 \\ As_1^n + Bs_2^n = 0 \end{array} \right\} \quad \begin{array}{l} A=-B \\ s_1^n - s_2^n = 0 \\ \downarrow \\ \cos(n\theta) + i \sin(n\theta) - \cos(n\theta) + i \sin(n\theta) = 0 \\ \downarrow \\ \sin(n\theta) = 0 \Rightarrow \theta = \frac{j\pi}{n} \quad (j=1 \dots n-1) \end{array}$$

So the possible eigenvalues are

$$\lambda_j = -\frac{1}{h} - s + 2 = 2 \left(1 - \cos \left(\frac{j\pi}{n} \right) \right) \quad j=1, \dots, n-1$$

$$\text{Recall: } \cos(2\alpha) = 1 - 2 \sin^2(\alpha) \Rightarrow 2 \sin^2(\alpha) = 1 - \cos(2\alpha)$$

$$\lambda_j = 4 \sin^2 \left(\frac{j\pi}{2n} \right) \quad j=1, \dots, n-1$$

$$\text{For the matrix } \mathbf{A} : \mu_j = \frac{1}{h^2} \lambda_j = \frac{4}{h^2} \sin^2 \left(\frac{j\pi}{2} \right) \quad \uparrow \quad m \cdot h = 1$$

Notice that the smallest eigenvalue is for $j=1$, and $\sin^2 \left(\frac{\pi}{2} \right) \approx \frac{h^2 \pi^2}{4}$

$$\boxed{\mu_{\min} \approx \frac{4}{h^2} \frac{h^2 \pi^2}{4} \approx \pi^2}$$

This proves the Proposition.

Now, we have some points to make:

$$(i) \quad \underline{e} = \mathbf{A}^{-1} \underline{u} \Rightarrow \|\underline{e}\|_2 \leq \|\mathbf{A}^{-1}\|_2 \|\underline{u}\|_2$$

(ii) \mathbf{A} is s.p.d. $\Rightarrow \mathbf{A}^{-1}$ is s.p.d. and $\|\mathbf{A}^{-1}\|_2$ is the spectral radius of \mathbf{A}^{-1}

$$(iii) \quad \text{Spectral radius of } \mathbf{A}^{-1} = \max \left(\left| \frac{1}{\text{eigen}(\mathbf{A})} \right| \right) = \frac{1}{\min \text{eigen}(\mathbf{A})} = \frac{1}{\mu_{\min}}$$

Proposition 1 +

(i)+(ii)+(iii) $\Rightarrow \|\mathbf{A}^{-1}\|$ is independent of h (stability)

The question now is: how $\|z\|_2$ scales with h ?

$$\|z\|_2^2 = \sum_{i,j} c_j^2 h^4 \leq \max_i (c_j^2 h^4) \sum_{i,j} 1 \leq \underbrace{\max_j c_j^2 h^4}_{\|z\|_\infty} \sum_{j=1}^n 1 = C_{\max} h^4 n = C_{\max} h^3$$

\uparrow
 $nh = 1$

$$\boxed{\|z\|_2 \sim O(h^{3/2})}$$

With the stability result we have: $\|\varepsilon\|_2 \sim O(h^{3/2})$

There are other ways to prove the stability and then the convergence. One is based on the definition of the Green function for the discrete problem (see Quarteroni - Sacco - Saleri)

So, one can prove that

$$\|\varepsilon\|_\infty \leq \frac{1}{8} \max_j |\varepsilon_j| \sim O(h^2)$$

There is another technique (similar to the Finite Element Method)

Let's define the discrete scalar product $(u, v)_h = h \sum_{j=0}^n c_j u_j v_j$

$c_{0,n} = \frac{1}{2}$
 $c_j = 1 \quad \text{if } 0 < j < n$

This is basically the **trapezoidal integration formula** of two functions with nodal values u_j, v_j

$$(T) \left(\int_0^l f g \right) = \frac{h}{2} (u_0 v_0 + u_1 v_1) + \frac{h}{2} (u_1 v_1 + u_2 v_2) + \dots + \frac{h}{2} (u_{n-1} v_{n-1} + u_n v_n).$$

Now, with some technicalities one can prove the following statements.

Denote

$$L_h u_h (v_h) = - \frac{u_{j+1} - 2u_j + u_{j-1}}{h^2}.$$

Then:

$$(1) \quad (L_h u_h, v_h)_h = (u_h, L_h v_h)_h \quad \forall u_h, v_h \text{ s.t. } u_0 = v_0 = u_n = v_n = 0$$

$$(2) \quad (L_h v_h, v_h)_h \geq 0 \quad \forall v_h \text{ s.t. } v_0 = v_n = 0$$

$$(3) \quad (L_h v_h, v_h)_h = \frac{1}{h} \sum_j (v_{j+1} - v_j)^2 = \boxed{\frac{h^2}{2} \sum_j \left(\frac{v_{j+1} - v_j}{h} \right)^2}$$

→ In general, this quantity is ≥ 0 . However, since we have $v_0 = v_n = 0$, we can say in detail that

$$(L_h v_h, v_h)_h = 0 \iff v_h = 0$$

Notice that $\frac{h^2}{2} \sum_j \left(\frac{v_{j+1} - v_j}{h} \right)^2$ is therefore a norm for the finite difference of order 1 of v_h (\approx first derivative)

$$(4) \quad (v_h, v_h)_h \leq \frac{1}{2} (L_h v_h, v_h)_h \quad \Rightarrow \quad \|v_h\|_h \leq \frac{1}{\sqrt{2}} \|\varepsilon_h\|_h$$

having set: $\|\varepsilon_h\|_h = (v_h, v_h)_h^{1/2}$ $\|\varepsilon_h\|_h = (L_h v_h, v_h)_h^{1/2}$

Now, if we put together the consequences of (1-4) we obtain:

Exact problem : $L_h u_h = f_j + c_j h^2$

Numerical problem $L_h u_h = f_j$

Error problem : $L_h e_h = c_j h^2$

$$(L_h e_h, e_h)_h = h^2 (c_j, e_h)_h = h^2 \left(h \sum c_j e_j \right) \leq h^2 h^2 \left(\sum c_j^2 \right)^{1/2} h^2 \left(\sum e_j^2 \right)^{1/2}$$

$$\|e_h\|_h^2 \leq h^2 C \|e_h\|_h$$

$$\sqrt{2} \|e_h\|_h^2 \leq h^2 C \|e_h\|_h$$

$$\Rightarrow \|e_h\|_h \leq \frac{C}{\sqrt{2}} h^2$$

Again we found that

$$\text{CONSISTENCY} + \text{STABILITY} \Rightarrow \text{CONVERGENCE}$$

In particular, in these linear problems, stability is equivalent to boudness. For the energy, we have

$$\sqrt{2} \|u_h\|_h^2 \leq \|u_h\|_h^2 \leq (L_h u_h, u_h)_h = (f, u_h)_h \leq \|f\|_h \|u_h\|_h$$

$$\|u_h\|_h \leq \frac{\|f\|_h}{\sqrt{2}}$$

This is the boudness that gives stability to the numerical problem.

Advection-Diffusion

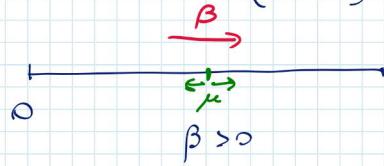
Problems

When we extend our method to more general problem, we can find that, in practice, stability and convergence are not enough. We are happy to know that the error vanishes with h , but certainly we do not want to take h too small (poor computer!!!)

Let us consider this problem :

$$\begin{cases} -\mu u'' + \beta u' = f & x \in (0,1) \\ u(0) = 0 \quad u(1) = 1 \end{cases}$$

This is a model problem for drift ($\beta u'$) + diffusion ($\mu u''$) problems.



μ = diffusion (Brownian motion)

β = drift/convection (or advection)

Model: pollutant in a river.

Before we study the numerical approximation, let's do some preliminary considerations.

Assume $\beta = 0$: $\begin{cases} -\mu u'' = 0 \\ u(0) = 0, \quad u(1) = 1 \end{cases}$ } Solution is rectilinear $u = x$

For $\mu = 0$, the solution is constant, but we have a problem with the boundary conditions.

For $\beta \neq 0$ we can compute the solution explicitly : Set $w = u'$ and assume $\beta \neq 0, \mu > 0$:

$$-\mu w' + \beta w = 0 \Rightarrow \frac{w'}{w} = \frac{\beta}{\mu} \Rightarrow \frac{d}{dx} \log(w) = \frac{\beta}{\mu}$$

$$\text{so } \log(w) = \frac{\beta}{\mu}x + A \Rightarrow w = e^{\frac{\beta}{\mu}x} B$$

↑
constant
↑
constant (e^A)

$$\text{Now } u' = w \Rightarrow u = B \frac{\mu}{\beta} e^{\frac{\beta}{\mu}x} + C$$

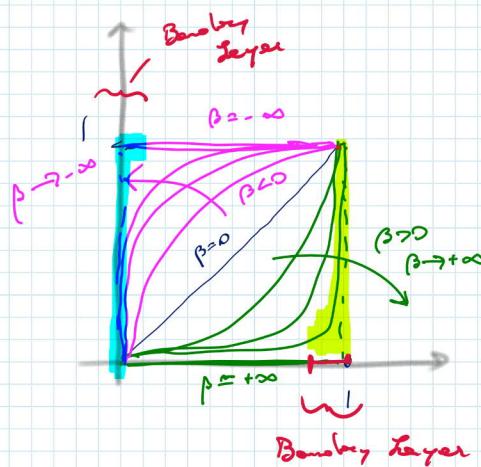
↑
constant
↑
constant

$$\left(\text{check: } u'' = B \frac{\beta}{\mu} e^{\frac{\beta}{\mu}x} \right)$$

$$u' = B e^{\frac{\beta}{\mu}x} \quad \left\{ -\mu u'' + \beta u' = 0 \right.$$

$$\left| \begin{array}{l} u(0) = 0 \\ u(1) = 1 \end{array} \right. \quad \left| \begin{array}{l} \downarrow \\ B \frac{\mu}{\beta} + C = 0 \\ \downarrow \\ B \frac{\mu}{\beta} e^{\frac{\beta}{\mu}} = 1 \end{array} \right.$$

$$\left. \begin{array}{l} \downarrow \\ B \frac{\mu}{\beta} (e^{\frac{\beta}{\mu}} - 1) = 1 \\ \downarrow \\ B = \frac{\mu}{\beta} (e^{\frac{\beta}{\mu}} - 1)^{-1} \end{array} \right.$$



In fluid mechanics, the part of the solution connecting the "almost" 0 (for $\beta > 0$) to the 1 (or the "almost" 0 for $\beta < 0$) is called BOUNDARY LAYER.

The length of the boundary layer scales with $\frac{1}{\beta}$. Larger is β vs μ , and smaller is the B.L.

We found the solution to the problem, and it is unique and bounded in O^{+1} .

Also, notice that

$$u' = \frac{\beta}{\mu} \frac{e^{\frac{\beta}{\mu}x}}{e^{\frac{\beta}{\mu}} - 1} \quad \text{so the derivative is bounded by } \frac{\beta}{\mu} \left(1 + \frac{1}{e^{\frac{\beta}{\mu}} - 1} \right)$$

How can we approximate this problem numerically?

Idea :

$$-\mu \frac{u_{j+1} - 2u_j + u_{j-1}}{h^2} + \beta \frac{u_{j+1} - u_{j-1}}{2h} = 0 \quad j = 1, \dots, m-1$$

$$u_0 = 0 \quad u_m = 1$$

Let's solve this explicitly, noting that each finite difference has a truncation error scaling with $O(h^2)$.

$$\left(\mu - \frac{\beta h}{2} \right) u_{j+1} - 2\mu u_j + \left(\mu + \frac{\beta h}{2} \right) u_{j-1} = 0 \quad (\text{assume } h : \mu - \beta h \neq 0)$$

$$\frac{\mu \pm \sqrt{\mu^2 - \mu^2 + \beta^2 \frac{h^2}{4}}}{\mu - \frac{\beta h}{2}} = \frac{1}{\frac{\mu + \frac{\beta h}{2}}{\mu - \frac{\beta h}{2}}}$$

$$u_j = A \varphi_1^j + B \varphi_2^j \quad u_0 = 0 \quad \Delta = -\beta$$

$$u_m = 1 \quad A = \left(\varphi_1^m - \varphi_2^m \right)^{-1}$$

Let's introduce the following adimensional number: $Peclet \equiv \frac{Re}{2\mu}$

Then, we can write: $\delta_1 = 1$, $\delta_2 = \frac{1+Re}{1-Re}$, so let

$$u_j = \frac{\delta_1^j - \delta_2^j}{\delta_1^n - \delta_2^n} = \frac{1 - \left(\frac{1+Re}{1-Re}\right)^j}{1 - \left(\frac{1+Re}{1-Re}\right)^n}$$

Notice that we have a numerical solution and it is possible to prove that the linear system

$$A \underline{u} = \underline{f}$$

with $A = \frac{\kappa}{h^2} \text{tridiag}([-1, 2, -1]) + \frac{\beta}{2h} \text{tridiag}([1, 0, -1])$ is invertible.

We will test that for $h \rightarrow 0$, the method has order 2 in $\|\cdot\|_\infty$ and $\|\cdot\|_1$.

However, we notice that for large h the solution features unphysical oscillations. Why?

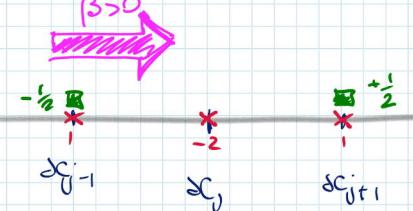
Note that for $Re > 1$ ($h > \frac{2\mu}{\beta}$), $\delta_2 < 0$, so δ_2^j is < 0 for j even and > 0 for j odd.

This introduces oscillations that make the solution unphysical.

Even if the solution is clear (take $h < \frac{2\mu}{\beta}$), this may be practically impossible:

in real problems we may have $\frac{h\mu}{\beta} \approx 10^{-6}$ or less!

This is not an instability of the numerical problem, yet for h large we need some practical solution. Even if not completely appropriate, these solutions go under the name of "stabilization".



$\times = 2^{\text{nd}}$ derivative

$\blacksquare = 1^{\text{st}}$ derivative

The reason of the "instability" is the lack of "physical consistency" of the scheme with the problem.

In this problem we have a "symmetric" dynamics (the diffusion) + a "directional" one (the convection).

The numerical scheme is "symmetric" and it does not include the directionality.

If you focus on x_j , the "wind" is blowing from the left to the right (for β positive), but the coefficients "upwind" (x_{j-1}) and "downwind" (x_{j+1}) are the same.

This is not "physically" consistent with the directional nature of the convective term.

The upwind Scheme

As a first idea, we break the "symmetry" by going upwind, i.e. in the direction where the wind is blowing from. For instance:

$$(\beta > 0): \quad u'(x_j) \approx \frac{u_j - u_{j-1}}{h}$$

$\Rightarrow \beta > 0$



\parallel We take data from "fresh" information

[When you go hunting, the prey must be "upwind" otherwise it smells you!]

The scheme reads:

$$-\frac{\mu}{h^2} u_{j+1} + \left(\frac{2\mu}{h^2} + \frac{\beta}{h} \right) u_j - \left(\frac{\mu}{h^2} + \frac{\beta}{h} \right) u_{j-1} = 0$$

$$-\frac{2\mu}{h^2} \frac{1}{2} u_{j+1} + \frac{2\mu}{h^2} \left(1 + \frac{\beta h}{2\mu} \right) u_j - \frac{2\mu}{h^2} \left(\frac{1}{2} + \frac{\beta h}{2\mu} \right) u_{j-1} = 0$$

$$(1 + \text{Pe}) \quad \left(\frac{1}{2} + \frac{\text{Pe}}{2} \right)$$

$$u_{j+1} - 2(1 + \text{Pe})u_j + (1 + 2\text{Pe})u_{j-1} = 0$$

$$\lambda_{1,2} = \frac{(1 + \text{Pe}) \pm \sqrt{1 + 2\text{Pe} + \text{Pe}^2 - 1 - 2\text{Pe}}}{2} =$$

$$= \begin{cases} 1 + 2\text{Pe} \\ 1 \end{cases} \quad \text{both } > 0$$

Test Upwind in the hands-on sessions.

The drawback of upwind is that the finite difference scheme we get is only 1st order:

We got stability by reducing accuracy (it happens sometimes)

In the numerical tests, we also notice that the numerical boundary layer is always overestimated vs. the exact one:



We want to explain why.

Notice that:

$$\beta \frac{u_j - u_{j-1}}{h} = \beta \frac{u_{j+1} - u_{j-1}}{2h} - \beta \frac{u_{j+1} - 2u_j + u_{j-1}}{2h} = \boxed{\beta \frac{u_{j+1} - u_{j-1}}{2h}} - \boxed{\frac{\beta h}{2} \frac{u_{j+1} - 2u_j + u_{j-1}}{h^2}}$$

Central FD extra numerical diffusion

UPWIND is like solving with central finite difference with an extra numerical diffusion: it is like solving the problem:

$$-\mu_h u'' + \beta u' = 0 \quad \text{with } \mu_h = \mu + \frac{\beta h}{2} = \mu (1 + \text{Pe})$$

Numerical Viscosity
(over diffusion)

The numerical BL is therefore $O\left(\frac{\mu_h}{\beta}\right) = O\left(\frac{\mu}{\beta} + \frac{h}{2}\right)$

From here, it is also evident that:

$$\text{Pe}_{\text{numerical}} = \frac{\beta h}{2\mu_h} = \frac{\beta h}{2\mu(1 + \text{Pe})} = \frac{\text{Pe}}{1 + \text{Pe}} < 1 \quad \forall h$$

Question: Can we buy stability without losing accuracy?

Two ingredients:

- (1) we need to go "upwind"
- (2) we need a 2nd order approximation of the first derivative.

Remark: Reaction problems

A complete real problem may include also a so called reactive term:

$$-\mu \frac{\partial^2 u}{\partial x^2} + \beta \frac{\partial u}{\partial x} + \sigma u = 0 \quad \begin{cases} x \in (0,1) \\ + B.C. \end{cases}$$

Diffusion Convection (Advection) Reaction

For instance, in the dynamics of a pollutant in a river, we have:

- diffusion is the dilution of the pollutant (random walk)
- convection is the movement of the water stream (β = water velocity)
- reaction is any chemical reaction with other agents (weighted by σ)

We found that if β is large vs μ we have numerical problems. What if $\sigma \gg \mu$?

For instance:

$$-\mu \frac{u_{j+1} - 2u_j + u_{j-1}}{h^2} + \sigma u_j = 0 \Rightarrow \mu u_{j+1} - (2\mu + h^2)u_j + \mu u_{j-1} = 0$$

$$\gamma^2 - 2\left(1 + \frac{\sigma h^2}{2\mu}\right)\gamma + 1 = 0 \quad \gamma_{1,2} = \frac{1 + \frac{\sigma h^2}{2\mu} \pm \sqrt{\left(\frac{\sigma h^2}{2\mu}\right)^2 - 1}}{2}$$

\downarrow
 No numerical oscillations!

$\gamma_{1,2} > 0 \quad \forall h$

(Reaction is less problematic than convection)

Other Types of Boundary Conditions

So far, we considered only conditions like $u(0) = \text{data}$. In fact, we may have different conditions.

Let's keep the example of the river pollutant.



$$\left. \begin{array}{ll} \text{Dirichlet (or Essential) Conditions:} & u(\text{boundary}) = \text{data} \\ \text{Neumann (or Natural) } \quad " \quad : & \mu \frac{\partial u}{\partial x}(\text{boundary}) = \text{data} \\ \text{Robin (or Mixed) } \quad " \quad : & \mu \frac{\partial u}{\partial x} + \chi u(\text{boundary}) = \\ & = \text{data} \end{array} \right\}$$

We prescribe a value for the concentration

We prescribe the flux of pollutant entering/leaving the domain

Typically $\text{data} = X \text{ Natural}$, so that:

$$\mu \frac{\partial u}{\partial x} = X \text{ Natural} - u$$

We prescribe a flux proportional to the difference between the external concentration X and the concentration at the boundary.

Let's see an example:

$$-\mu \frac{\partial^2 u}{\partial x^2} + \beta \frac{\partial u}{\partial x} + \sigma u = 0$$

$$u(0) = 1$$

$$\left[\mu \frac{\partial u}{\partial x} + \chi u \right](1) = X u_{\text{ex}} \quad (\text{for } \chi = 0 \text{ we have Neumann homogeneous})$$

$$\frac{\partial u}{\partial x}(1) + \chi u(1) = X u_{\text{ex}}$$



A first order (stable) scheme :

Assuming $\beta > 0$:

$$\left\{ \begin{array}{l} -\frac{u_{j+1} - 2u_j + u_{j-1}}{h^2} + \beta \frac{u_j - u_{j-1}}{h} + \sigma u_j = 0 \quad j = 1, \dots, n-1 \\ u_0 = 1 \\ -\mu \frac{u_{n-1} - u_n}{h} + \chi u_n = \chi u_{\text{exact}} \end{array} \right.$$

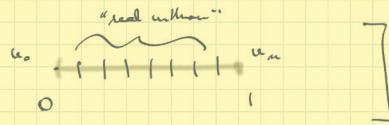
We can set a $(n+1) \times (n+1)$ system
solving this problem.

What if we want a 2nd order stable scheme? (HW2)

REMARK

The Dirichlet conditions $u(\text{boundary}) = \text{data}$ are quite immediate to consider: we don't have an unknown on the boundary!!!

When we work on paper, we simply "eliminate" these unknowns and work on a smaller system:



We eliminate u_0 and u_m and solve an $(n-1) \times (n-1)$ system.

In practice, in real coding, we prefer another strategy: we retain

$$\begin{cases} u_0 = u_{\text{left}} \\ u_1 = u_{\text{right}} \end{cases}$$

as equations and solve a $(n+1) \times (n+1)$ system

In a small example:

$$\left\{ \begin{array}{l} \alpha_{11}x_1 + \alpha_{12}x_2 + \alpha_{13}x_3 = f_1 \\ \alpha_{21}x_1 + \alpha_{22}x_2 + \alpha_{23}x_3 = f_2 \\ \alpha_{31}x_1 + \alpha_{32}x_2 + \alpha_{33}x_3 = f_3 \end{array} \right.$$

with $\alpha_1 = \alpha$, $\alpha_2 = \beta$

$$x_{22}x_2 = f_2 - \alpha_{21}\alpha - \alpha_{23}\beta \quad (\text{for } \alpha_{22} \neq 0)$$

Theory

Practice

$$\begin{aligned} 8x_1 &= f_2 \\ \alpha_{21}x_1 + \alpha_{22}x_2 + \alpha_{23}x_3 &= f_2 \\ 8x_3 &= \beta \end{aligned}$$

This allows a better management of the memory and coding.

Solving a system $(n+1) \times (n+1)$ has basically the same cost as $(n-1) \times (n-1)$.

This coefficient intend to be of the same order of $\alpha_{21}, \alpha_{22}, \alpha_{23}$
so to have good conditioning properties of the matrix.

In general, when ^{a PDE problem}, we have four steps:

Pre-processing : Definition of physical and numerical parameters

Assembly : (1) Construction of the matrix A and the right hand side b w/out the boundary conditions

(2) Prescription of the boundary conditions

→ The two-step structure allows a more effective coding (no "if" statements in "for" loops).

Solving : the solution of the linear system(s)

Post-processing : Visualization / Verify

Week 3

Parabolic Problems: Heat (Diffusion) Equation and Similar

Lamgautier - Chap. 3
Quarteroni - Chap. 5

INTRODUCTION

In this Chapter, we consider the FD approximation of the famous equation:

$$\frac{\partial u}{\partial t} - \mu \frac{\partial^2 u}{\partial x^2} = f(x,t) \quad x \in (0,1), t > 0 \\ (\mu > 0 \text{ is a constant})$$

with Boundary Conditions:

$$u(0) = u(1) = 0 \quad \forall t > 0$$

or Initial Condition:

$$u(x,0) = u_0(x).$$

It is, basically, the unsteady version of the problem considered in Week 2.

We do not consider the analytical solution, it can be obtained by separation of variables (see undergraduate courses on PDEs).

However, we give a little insight to the analysis of the solution, even without finding it explicitly. We will find a similar pattern in the analysis of the numerical solution.

Then, we will recall some basics for finding the analytical solution.

Physical Meaning of the equation

This is the classical "diffusion" equation: it may represent the evolution of a pollutant in time, driven only by the random molecular agitation. We can generalize it to the unsteady Advection-Diffusion-Reaction:

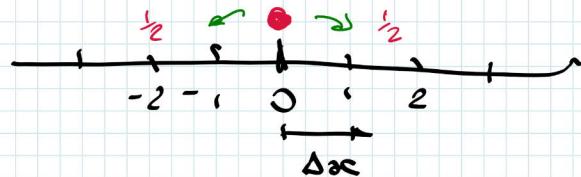
$$\frac{\partial u}{\partial t} - \mu \frac{\partial^2 u}{\partial x^2} + \beta \frac{\partial u}{\partial x} + \sigma u = f(x,t), \quad x \in (0,1), t > 0 \\ + \text{B.C.} + \text{I.C.} \\ u(0,t) = u(1,t) = 0 \quad u(x,0) = u_0(x)$$

Another interpretation is the "heat equation": in this case, $u(x,t)$ is the temperature of a rod. The temperature evolves for the thermal conductivity, represented by the coefficient μ . In this case, the boundary conditions have the following meaning:

- (1) $u(\text{boundary}) = \text{given} \Rightarrow \text{we measure the temperature on the boundary}$
- (2) $\mu \frac{\partial u}{\partial x}(\text{boundary}) = \text{given} \Rightarrow \text{we prescribe the thermal flux} (= 0 \text{ for insulated boundary})$
- (3) $\mu \frac{\partial u}{\partial x} + \chi u = \gamma u_{\text{ext}} \Rightarrow \text{Radiation conditions.}$
(thermal flux proportional to the difference between internal and external temperature).

It is worth recalling the following fact.

If you consider a 1D random walk, i.e. a game where a mass is located at the center of an interval (say in 0) and at each time step the mass can move to the left with probability q and to the right with probability $1-q$.



Denoting by Δx the step between two cells or positions and by Δt the time step, and by $p(x_j, t^n)$ the probability that the mass is in $x_j = j \Delta x$ at time $t^n = n \Delta t$, for $\Delta t \rightarrow 0$ and $\Delta x \rightarrow 0$ this probability obeys the diffusion equation with the initial condition

$$p(x, 0) = \begin{cases} 1 & x=0 \\ 0 & x \neq 0 \end{cases}$$

This object is not a function, it is a distribution (called Dirac- δ).

This justifies the attention in the Langtangen book to 1D Random Walks in Chapter 3. The fundamental solution of the heat equation on $x \in \mathbb{R}$ (no boundary conditions) with the Dirac- δ as initial condition is a Gaussian function - not surprisingly.

For a complete proof, see S. Salsa, PDEs in Action, Chap. 2

To prove this, you need to postulate that:

$$\lim_{\substack{\Delta t \rightarrow 0 \\ \Delta x \rightarrow 0}} \frac{\Delta x^2}{\Delta t} = \text{constant}$$

and this constant turns out to be μ .

Analyzing the solution w/out computing it

Let's assume we know that a solution exists. What can we say about it, without computing it explicitly? The techniques we see here will be useful to understand also one numerical solution.

We know that u is such that

$$\frac{\partial u}{\partial t} - \mu \frac{\partial^2 u}{\partial x^2} = f, \quad u(0) = u(1) = 0 \quad u(x, 0) = u_0(x).$$

For simplicity, let's start with the case $f = 0$.

If we multiply the equation by u , we integrate over $x \in (0, 1)$ and we integrate by parts, we have:

$$\int_0^1 u \frac{\partial u}{\partial t} - \mu \int_0^1 u \frac{\partial^2 u}{\partial x^2} = 0 \Rightarrow \frac{1}{2} \int_0^1 \frac{\partial(u^2)}{\partial t} + \mu \int_0^1 \left(\frac{\partial u}{\partial x} \right)^2 = 0$$

In fact:

$$\begin{aligned} u \frac{\partial u}{\partial t} &= \frac{1}{2} \frac{\partial(u^2)}{\partial t} \\ -\mu \int_0^1 u \frac{\partial^2 u}{\partial x^2} &= -\mu \left[\frac{\partial u}{\partial x} \right]_0^1 + \mu \int_0^1 \frac{\partial u \partial u}{\partial x \partial x} \end{aligned} \Rightarrow \frac{1}{2} \frac{\partial}{\partial t} \int_0^1 u^2 + \mu \int_0^1 \left(\frac{\partial u}{\partial x} \right)^2 = 0 \quad *$$

Using a notation that will be clear better with Finite Elements, for a generic function g (regular enough):

$$\|g\|^2 \stackrel{\text{def}}{=} \int_0^1 g^2 dx$$

Then \star needs : $\frac{d}{dt} \|u\|^2 + \mu \left\| \frac{\partial u}{\partial x} \right\|^2 = 0$

Now, let's integrate in time between 0 and a final time T_f (generic)

$$\int_0^{T_f} \frac{d}{dt} \|u\|^2 dt + \int_0^{T_f} \mu \left\| \frac{\partial u}{\partial x} \right\|^2 dt = 0$$

$\underbrace{\quad}_{\geq 0}$

$$\|u\|^2(T_f) + \int_0^{T_f} \mu \left\| \frac{\partial u}{\partial x} \right\|^2 dt = \|u\|^2(0) = \|u_0\|^2$$

$$\Rightarrow \|u\|(T_f) \leq \|u_0\|$$

↑
given

The solution is bounded by the initial condition!

If we have $f=0$, we have the following inequality, following similar steps :

$$\|u\|^2(T_f) + \int_0^{T_f} \mu \left\| \frac{\partial u}{\partial x} \right\|^2 dt = \|u_0\|^2 + \int_0^{T_f} \int_0^1 f u dx dt$$

The Cauchy-Schwarz inequality states that : $\left| \int_0^1 f u \right| \leq \|f\| \|u\|$

So we have :

$$\|u\|^2(T_f) \leq \|u_0\|^2 + \int_0^{T_f} \|f\| \|u\| dt$$

At this point we call a result called : Gronwall Lemma :

GRONWALL LEMMA

If $y(t) \leq f(t) + \int_0^t g(s)y(s) ds$ } then $y(t) \leq f(t) + \int_0^t f(s) e^{\int_s^t g(s) ds} ds$
 with $f \geq 0, g \geq 0$

By a direct application of this lemma, we have that the solution is bounded.

Notice that there exists a "discrete" version of this lemma (we can use it for our numerical solution).

Calculating the exact solution (in 1D)

In 1D we can compute the exact solution using the method of SEPARATION OF VARIABLES.

We recall here the basic principles.

(1) We use the superposition of effects (that holds because the problem is linear) to postulate a solution in the form :

$$u(x,t) = \sum_{j=1}^{\infty} X_j(x) T_j(t)$$

(2) We work out each component .

Let's look for the component $X_j(x) T_j(t)$ and assume that we can also write $f(x,t) = \sum F_j(t) X_j(x)$. We will see in a second who are the $X_j(x)$.

If we plug our solution in the equation, we have:

$$(Δ) \quad \sum_{j=1}^{\infty} \frac{d\bar{T}_j}{dt} X_j - \mu \sum_{j=1}^{\infty} \frac{d^2 X_j}{dx^2} \bar{T}_j = \sum_{j=1}^{\infty} F_j X_j$$

with $u_0(x) = \sum_{j=1}^{\infty} C_j X_j(x)$

Now, let's assume we can solve the problem:

$$(□) \quad \frac{d^2 X_j}{dx^2} = -\lambda_j X_j \quad \text{with } X_j(0) = X_j(1) = 0$$

Then we have from (□)

$$\sum_{j=1}^{\infty} \left(\frac{d\bar{T}_j}{dt} + \mu \lambda_j \bar{T}_j - \bar{F}_j \right) X_j = 0$$

So that we solve the problem by solving the ODEs:

$$\frac{d\bar{T}_j}{dt} = -\mu \lambda_j \bar{T}_j + \bar{F}_j \quad \text{with } \bar{T}_j(0) = C_j$$

Can we solve (□)?

The specific answer is easy: non trivial X_j are in the form:

$$X_j = \sin(\sqrt{\lambda_j} x) \quad \text{with } \sqrt{\lambda_j} = j\pi$$

In general the problem (□) is well understood and it goes under the name of STURM-LIOUVILLE EIGENVALUE PROBLEM.

It is possible to prove that the eigenfunctions X_j form an orthogonal basis, i.e., for a generic function $g(x)$ with enough regularity, we can always write:

$$g(x) = \sum_{j=1}^{\infty} G_j X_j(x) \quad \text{with } \int_0^1 X_j X_k = 0 \quad \text{for } j \neq k$$

and $G_j = \frac{\int_0^1 g(x) X_j(x)}{\int_0^1 X_j^2}$ (generalized Fourier expansion).

So to summarize:

$$(1) \text{ one computes: } X_j : \frac{d^2 X_j}{dx^2} = -\lambda_j X_j + \text{B.C.}$$

$$(2) \text{ Then: } C_j = \frac{\int_0^1 u_0(x) X_j(x)}{\int_0^1 X_j^2}, \quad \bar{F}_j = \frac{\int_0^1 f(x) X_j(x)}{\int_0^1 X_j^2}$$

$$(3) \text{ Then: } \frac{d\bar{T}_j}{dt} = -\mu \lambda_j \bar{T}_j + \bar{F}_j \quad T_j(0) = C_j$$

$$(4) \quad u(x,t) = \sum_{j=1}^{\infty} \bar{T}_j X_j$$

Finite Differences for the Heat Equation

The heat equation is the first truly PDE problem we have to solve : we take advantage of the experience of the previous two chapters.

$$\frac{\partial u}{\partial t} - \mu \frac{\partial^2 u}{\partial x^2} = f$$

We need to discretize both the space and the time derivative.

SPACE (SEMI) DISCRETIZATION

Let $u_j(t)$ be the approximation of $u(x_j, t)$ and let us collocate the problem in x_j , with the usual approximation for the 2nd derivative :

$$(•) \quad \frac{du_i}{dt} - \frac{\mu}{h^2} (u_{j+1}(t) - 2u_j(t) + u_{j-1}(t)) = f(x_j, t)$$

with $u_j(t=0) = u_0(x_j)$ and $u_0(t) = u_n(t) = 0$.

Notation
 $h = \Delta x$

If we define the vector (function of time) :

$$\underline{u}(t) = \begin{bmatrix} u_0(t) \\ u_1(t) \\ \vdots \\ u_n(t) \end{bmatrix}$$

we can recognize in the semi-discrete problem (•) the system of ODEs :

$$\frac{d\underline{u}}{dt} = A\underline{u} + \underline{f} \quad (\spadesuit) \quad (\text{with } u_j(0) = u_0(x_j))$$

where :

$$A = \frac{\mu}{h^2} \text{tridiag}(1, -2, 1) \quad \underline{f} = \begin{bmatrix} f(x_0, t) \\ \vdots \\ f(x_n, t) \end{bmatrix}$$

We know that $\text{tridiag}(-1, 2, -1)$ is s.p.d.

Since $\frac{\mu}{h^2} > 0$, we conclude that A is s.m.d. (negative) and we already know that the eigenvalues are all real and negative

(if a generic matrix B has eigenvalues $1, -3$ has eigenvalues -1)

How do we discretize in time the system (•) ?

The answer is an entire chapter of MATH 516 !

We faced a 2nd order Cauchy problem in Week 1. Here we have a first order problem and the solution is simpler.

It is particularly "educational" to look at a simple class of methods called Θ -methods. In fact, Runge-Kutta methods are always an excellent option.

Θ -Methods

The following presentation is different if compared with the books, but I think it is more rigorous.

The integral form of the Cauchy problem reads:

$$\int_{t_1}^{t_2} \frac{du}{dt} = \int_{t_1}^{t_2} A u + f \cdot \Rightarrow u(t_2) = u(t_1) + \int_{t_1}^{t_2} (A u + f) dt$$

Two words to recall bases of numerical quadratures (integrals).

We can approximate the integral of a generic function $y(f)$ as follows:

$$\int_{t_1}^{t_2} y(t) dt \approx w_1 y(a) + w_2 y(b)$$

with w_1, w_2 "weights" and $t_1 \leq a \leq b \leq t_2$.

If w_1, w_2, a and b can be all be determined we can maximize the accuracy of the formula (Gauss).

Also, notice that in order to find the correct integral of a constant,

$$w_1 + w_2 = (t_2 - t_1)$$

In fact, if y is constant:

$$\int_{t_1}^{t_2} y = \text{constant} \times (t_2 - t_1) = (w_1 + w_2) \times \text{constant}$$

So, from now on we set:

$$w_1 = (1-\theta)(t_2 - t_1) \quad w_2 = \theta(t_2 - t_1).$$

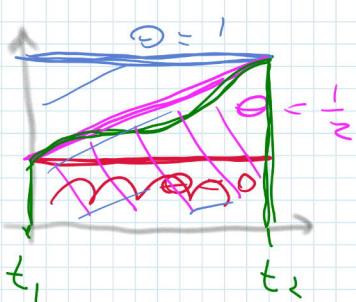
The nodes a and b can be selected "optimally" (for the accuracy) or with any practical convenience.

A common choice is: $a = t_1$ $b = t_2$.

Then, we have:

$$\int_{t_1}^{t_2} y \approx (y(t_1)(1-\theta) + y(t_2)\theta)(t_2 - t_1)$$

For $\theta = 0, 1$ we approximate the integral with a rectangular rule (in red in the figure), for $\theta = \frac{1}{2}$ we have the trapezoidal rule (in magenta).



Notice that the only formula exact also for linear functions is for $\Theta = \frac{1}{2}$

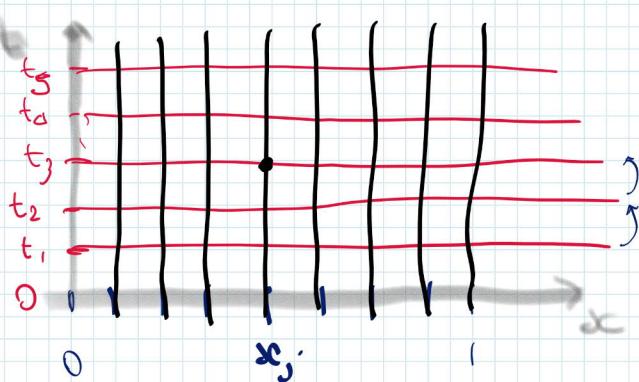
$$\frac{1}{2} = \int_0^1 x e^{xt} dx = \Theta \cdot \Theta + (1-\Theta) \cdot 1 = 1 - \Theta$$

$$\boxed{\Theta = \frac{1}{2}}$$

The best accuracy is attained with $\Theta = \frac{1}{2}$.

Full Discretization

Now, let's discretize the time axis and apply the Θ method to solve at each time step:



$$\begin{aligned} t^{n+1} &= t^n + \Delta t \\ \frac{du}{dt} &= A u + f \\ u^{n+1} &= u^n + \int_{t^n}^{t^{n+1}} (A u + f) dt \approx \\ &\approx u^n + \Delta t A u^n (1-\Theta) + \Delta t f^n (1-\Theta) \\ &\quad + \Delta t A u^{n+1} \Theta + \Delta t f^{n+1} \Theta \\ \text{where } f^n &= \left[\begin{array}{c} f(x_0, t^n) \\ f(x_1, t^n) \\ \vdots \\ f(x_{m-1}, t^n) \end{array} \right] \end{aligned}$$

We can rewrite the step as:

$$(I - \Delta t \Theta A) u^{n+1} = (I + \Delta t (1-\Theta) A) u^n + \Delta t (\Theta f^{n+1} + (1-\Theta) f^n)$$

Notice that for $\Theta = 0$ we do not have a linear system to solve.
for $\Theta \neq 0$ we need to solve a linear system **AT EACH TIME STEP**:
it is critical to have an efficient linear solver.

On the choice of the linear solver ($\Theta \neq 0$)

The system has matrix:

$$I - \Delta t \Theta A = \text{tridiag} \left[1, 1 - 2\frac{\Delta t}{h^2} \Theta, 1 \right]$$

It's a tridiagonal matrix: sparse with structured pattern
It's a time-independent matrix

If we perform an LU factorization, this holds for all the time step!
So, a reasonable choice for the linear solver is:

$$(1) \text{ Factorize } (I - \Delta t \Theta A) = LU$$

(2) Time Loop :

for $t_j < T_f$:

$$I + U \underbrace{U}_{\sim}^{\sim \sim \sim \sim} = (I + \Delta t (1-\Theta) A) \underbrace{U}_{\sim}^{\sim} + \underbrace{\Theta f}_{\sim}^{\sim \sim \sim}$$

two easy systems (Order of \sim operations each)

}

The real questions now are :

(1) is the method stable?

(2) how the accuracy changes with Θ , A and Δt ?

Stability Analysis of the Implicit Euler Scheme

Let's preliminary consider the Implicit / Backward Euler scheme : we will see that it is unconditionally stable. The technique resembles the one used for the analytical solution.

$$\underline{u}^n - \Delta t A \underline{u}^n = \underline{u}^{n-1} + \underline{f}^n$$

(1) Multiply the two sides by \underline{u}^n

$$\|\underline{u}^n\|_2^2 - \Delta t \underline{u}^n A \underline{u}^n = \underline{u}^n \cdot \underline{u}^{n-1} + \underline{u}^n \cdot \underline{f}^n$$

(2) Recall that $-A$ is s.p.d.

$$\Rightarrow -\Delta t \underline{u}^n A \underline{u}^n > 0$$

(3) Also, recall that

$$0 \leq \|\underline{u}^n - \underline{u}^{n-1}\|_2^2 = (\underline{u}^n - \underline{u}^{n-1})^T (\underline{u}^n - \underline{u}^{n-1}) = \|\underline{u}^n\|_2^2 + \|\underline{u}^{n-1}\|_2^2 - 2 \underline{u}^n \cdot \underline{u}^{n-1}$$

$$\underline{u}^n \cdot \underline{u}^{n-1} \leq \frac{1}{2} \|\underline{u}^n\|_2^2 + \|\underline{u}^{n-1}\|_2^2$$

$$\Rightarrow \frac{1}{2} \|\underline{u}^n\|_2^2 - \Delta t \underline{u}^n A \underline{u}^n \leq \frac{1}{2} \|\underline{u}^{n-1}\|_2^2 + \underline{u}^n \cdot \underline{f}^n$$

Now, we sum for $n = 1, 2, \dots, N_F$ where $N_F = \frac{T_F}{\Delta t}$

and we get

$$\|\underline{u}^{N_F}\|^2 + \|\underline{u}^{N_F-1}\|^2 + \dots + \|\underline{u}^1\|_2^2 - 2\Delta t \sum_{n=1}^{N_F} \underline{u}^{n-1} A \underline{u}^n \leq \|\underline{u}^{N_F-1}\|^2 + \|\underline{u}^{N_F-2}\|^2 + \dots + \|\underline{u}^0\|^2 + 2 \sum_n \|\underline{f}^n\| \|\underline{u}^n\|$$

(Discrete)
Cauchy-Schwarzs inequality

After cancelling the corresponding terms, and dropping the positive term:

$$\|\underline{u}^*\|^2 \leq \|\underline{u}^0\|^2 + 2 \sum_n \|\underline{u}^n\| \|\underline{f}^n\|$$

\Rightarrow For $\underline{f} = 0$ we have a bound on the solution at any time, based on initial date:

THE SOLUTION CANNOT EXPLODE!

\Rightarrow For $\underline{f} \neq 0$ we get to the bound by using the "discrete" equivalent of the Gronwall Lemma

\Rightarrow Backward Euler is unconditionally stable.

Stability Analysis of the Θ -Methods

The technique for Backward Euler works because the term

$$-2\Delta t \sum_{n=1}^{N_F} \underline{u}^{n-1} A \underline{u}^n > 0$$

What can we do in general?

Some BASIC FACTS

If $A = \Delta t \frac{\mu}{h^2}$ tridiag $(1, -2, 1)$ is s. p. negative:

specifically we found that it has all negative and distinct eigenvalues: $\lambda_j = -\frac{\Delta t \mu}{h^2} 4 \sin\left(jh \frac{\pi}{2}\right)$ $j = 1, \dots, N_x - 1$

where N_x is the number of space intervals ($N_x \cdot h = 1$).

When a matrix has all distinct eigenvalues, then it is diagonalizable, i.e., there exists a matrix T s.t.:

$$AT = TD \quad D = \text{diag}(\lambda_j)$$

T is invertible, and in the case of a symmetric matrix, it is also orthogonal, so we can write:

$$T^{-1}AT = D \Rightarrow T^TAT = D$$

because

$$TT^T = I \Rightarrow T^T = T^{-1}.$$

We are ready to analyze the general Θ method

(The analysis of Saenger book is different: Von Neumann)

Generic step of the Θ Method:

$$\textcircled{1} \quad (\mathbf{I} - \Delta t \Theta A) \underline{u}^n = (\mathbf{I} + \Delta t (1-\Theta) A) \underline{u}^{n-1} + \Theta \underline{f}^n + (1-\Theta) \underline{f}^{n-1}$$

Let's target the asymptotic behavior of the solution when $\underline{f} = 0$. The exact solution is a dissipative one, so we know that for $\underline{f} = 0$, $\underline{u}_{\infty} \rightarrow 0$.

In fact, for $\underline{f} = 0$:

$$\frac{d\underline{u}}{dt} = A\underline{u} \quad \text{with eigenvalues of } A \text{ real and negative.}$$

Now, let's rewrite the step $\textcircled{1}$ as follows:

$$(\mathbf{T}^{-1} - \Delta t \Theta T^{-1} \boxed{T^{-1} A T T^{-1}}) \underline{u}^n = (\mathbf{T}^{-1} + \Delta t (1-\Theta) T \boxed{T^{-1} A T T^{-1}}) \underline{u}^{n-1}$$

$$\Downarrow$$

$$T(\mathbf{I} - \Delta t \Theta D) T^{-1} \underline{u}^n = T(\mathbf{I} + \Delta t (1-\Theta) D) T^{-1} \underline{u}^{n-1}$$

T is invertible; also, denote:

$$T^{-1} \underline{u}^n = \underline{w}^n \quad \text{and} \quad T^{-1} \underline{u}^{n-1} = \underline{w}^{n-1}$$

Notice that if we know \underline{w}^j ($\forall j$), then we know \underline{v}^j
 $(\underline{v}^j = \underline{T} \underline{w}^j)$
 and if \underline{w}^j blows up, so does \underline{v}^j (and the other way round)

so we have:

$$(\mathbf{I} - \Delta t \Theta D) \underline{w}^{\sim} = (\mathbf{I} + \Delta t(1-\Theta) D) \underline{w}^{\sim-1}$$

But this is a diagonal system, so we can write the single equation:

$$(1 - \Delta t \Theta \lambda_j) w_j^{\sim} = (1 + \Delta t(1-\Theta) \lambda_j) w_j^{\sim-1}$$

$$w_j^{\sim} = \frac{1 + \Delta t(1-\Theta) \lambda_j}{1 - \Delta t \Theta \lambda_j} w_j^{\sim-1} \Rightarrow$$

$$\Rightarrow w_j^{\sim} = \left(\frac{1 + \Delta t(1-\Theta) \lambda_j}{1 - \Delta t \Theta \lambda_j} \right)^n w_j^{\circ}$$

We have that the solution remains stable if and only if

$$\left| \frac{1 + \Delta t(1-\Theta) \lambda_j}{1 - \Delta t \Theta \lambda_j} \right| < 1$$

recall that
 $\lambda_j < 0$

\Downarrow

$$-1 + \Delta t \Theta \lambda_j < 1 + \Delta t(1-\Theta) \lambda_j < 1 - \Delta t \Theta \lambda_j$$

$$\underbrace{|1 - \Delta t(1-\Theta)\lambda_j|}_{1 - \Delta t(1-\Theta)|\lambda_j|} < |1 + \Delta t\Theta\lambda_j|$$

$$|\Delta t\lambda_j| + \Theta\Delta t|\lambda_j| < 1 + \Delta t\Theta|\lambda_j|$$

ALWAYS TRUE

$$-\Delta t\Theta|\lambda_j| + \Delta t(1-\Theta)|\lambda_j| < 2 \quad \text{if } \Theta > \frac{1}{2} \text{ always true}$$

$$\Delta t|\lambda_j|(1-2\Theta) < 2 \Rightarrow$$

if $\Theta < \frac{1}{2}$ true for $\Delta t < \frac{2}{|\lambda_j|(1-2\Theta)}$

Conclusion 1 (Theorem)

The Θ method for the heat equation is

(1) unconditionally absolutely stable for $\Theta > \frac{1}{2}$

(2) conditionally stable for $\Theta < \frac{1}{2}$ with the condition:

$$\Delta t \leq \frac{1}{|\lambda_j|} \frac{2}{1-2\Theta} \quad j=1, \dots, N_x$$

This condition must be taken in the most restrictive sense:

$$\Delta t \leq \min_j \left| \frac{1}{\lambda_j} \right| \frac{2}{1-2\Theta} = \frac{1}{\max_j |\lambda_j|} \frac{2}{1-2\Theta}$$

Now, recalling that $|\lambda_j| \sim \frac{\alpha}{h^2} \sin^2\left(\frac{j\pi}{2} h\right)$ we conclude that $\max_j \sin^2\left(\frac{j\pi}{2} h\right) \underset{j=N_x-1}{\approx} \sin^2\left(\frac{\pi}{2}(1-\epsilon)\right) \approx 1 - O(\epsilon^2)$ so that

$$\max_j |\lambda_j| \sim \frac{1}{h^2}$$

Conclusion 2

For $\Theta < \frac{1}{2}$, the stability condition relates the discretization parameters:

$$\Delta t \leq C Q^2$$

Overall, this is a local constraint: if you want to refine your spatial mesh ($h \rightarrow \frac{h}{2}$), the temporal mesh must be refined more ($\Delta t \rightarrow \frac{\Delta t}{4}$).

For this reason, methods with $\Theta \geq \frac{1}{2}$ are more popular.

We will check this conclusions with the hands-on session.

Accuracy Analysis.

We just give some results with a "sketch" of the proof.

The space discretization has an error $O(h^2)$

The time discretization has order (we didn't prove this here):

$$O(\Delta t^p) \text{ where } p \text{ depends on } \Theta : \begin{cases} P\left(\frac{1}{2}\right) = 2 \\ P(\Theta) = 1 \quad \forall \Theta \neq \frac{1}{2} \end{cases}$$

We can prove the theorem:

Θ methods with centred discretization of the 2nd derivative have an error:

$$O\left(\Delta t^{P(\Theta)} + h^2\right) \text{ with } \begin{cases} P\left(\frac{1}{2}\right) = 2 \\ P(\Theta) = 1 \quad \forall \Theta \neq \frac{1}{2} \end{cases}$$

As usual, we get to the proof by combining stability and consistency.

Remarks

(1) We have two sources of errors, one controlled by Δt , the other by h .

In the numerical experiments it can be not easy to test the order.

You can:

(a) change Δt and h simultaneously
 $(\Delta t \rightarrow \frac{\Delta t}{2} \text{ and } h \rightarrow \frac{h}{2})$

(b) take an extremely small parameter (e.g., Δt) and play with the other (h) so to "kill" one error and play with the other

(2) The errors of the time discretization is generally of "dissipative" nature (in particular for Implicit Euler): they tend to dissipate the energy of the problem.

Keep in mind this in practice sessions!!

A Final Word on More Complex Problems.

We may now consider the problem:

$$\frac{\partial u}{\partial t} - \mu \frac{\partial^2 u}{\partial x^2} + \beta \frac{\partial u}{\partial x} + \sigma u = f \quad x \in (0,1), t > 0$$

+ B.C. + I.C.

In this problem, u may be the concentration of pollutant in a river:

μ is the diffusivity of the pollutant

β the velocity of the river

& the reaction of the pollutant with other agents.

A possible set of boundary conditions is

$$u(0, t) = g(t)$$

$$\frac{\partial u}{\partial x}(1, t) = 0$$

(the concentration is space invariant at the outlet)

Then, one can replace the condition in 1 with:

$$\frac{1}{h} \left(\frac{3}{2} u_{N_x} - 2 u_{N_x-1} + \frac{1}{2} u_{N_x-2} \right) = 0$$

and write the system:

$$\left\{ \begin{array}{l} u_0^n = g(t^n) \\ u_j^n = u_j^{n-1} - \Delta t \theta \left(\frac{\mu}{h^2} (u_{j+1}^n - 2u_j^n + u_{j-1}^n) + \frac{\beta}{2h} (u_{j+1}^n - u_{j-1}^n) + \sigma u_j^n \right) \\ \quad - \Delta t (1-\theta) \left(\frac{\mu}{h^2} (u_{j+1}^{n-1} - 2u_j^{n-1} + u_{j-1}^{n-1}) + \frac{\beta}{2h} (u_{j+1}^{n-1} - u_{j-1}^{n-1}) + \sigma u_j^{n-1} \right) \\ \quad + \Delta t \theta f_j^n + \Delta t (1-\theta) f_j^{n-1} \\ 2 u_{N_x} - 2 u_{N_x-1} + \frac{1}{2} u_{N_x-2} = 0 \end{array} \right.$$

for $j = 1, \dots, N_x-1$

This is expected to be a scheme of order $O(h^2 + \Delta t^{P(\Theta)})$ unconditionally stable for $\Theta \geq \frac{1}{2}$.

Remark

Of course, the centred scheme may suffer from oscillations when $\text{Pe} > 1$, so, either we enforce $\text{Pe} < 1$ with h small enough, or we need some upwind scheme (possibly preserving the accuracy).

However, notice that: $\frac{\partial u}{\partial t}(t^n) \approx \frac{u^n - u^{n-1}}{\Delta t} - \underbrace{\frac{\Delta t \partial^2 u}{2 \Delta t^2}}_{\text{H.O.T.}}$

dropping this term amounts to adding a form of numerical dissipation that can delay or even suppress the Pe-related instabilities.

However, it is safe to enforce $\text{Pe} < 1$.

Week 4

Hyperbolic Problems

(1) Scalar Conservation Laws

(Linear)

for non-linear problems, see specific books)

(2) Systems of Conservation Laws

(3) Wave Equation

Most of the material is in the slides.

Lagrange book: Chapt. 4 and 2

Quarteroni: Chapt. 13

Quarteroni-Valli: Chapt. 14

Here I address specific topics not covered in the slides

(1) Derivation of the Lax-Wendroff Scheme

The LW scheme is an explicit scheme very popular in this field.

The derivation of the scheme has a STRONG EDUCATIONAL CONTENT that goes beyond the scheme itself.

Equation:

$$\frac{\partial u}{\partial t} + \alpha \frac{\partial u}{\partial x} = 0 \quad x \in \mathbb{R}, \quad t > 0$$

LW is an example of "Taylor-based" method.

Taylor expansion of the solution in time:

$$u(x, t^{n+1}) = u(x, t^n) + \frac{\partial u}{\partial t}(x, t^n) \Delta t + \frac{\partial^2 u}{\partial t^2}(x, t^n) \frac{\Delta t^2}{2} + \text{H.O.T.}$$

This is true for any regular function.

But in the case of our problem:

$$\frac{\partial u}{\partial t} = -\alpha \frac{\partial u}{\partial x} \quad (\text{from the equation})$$

and

$$\begin{aligned} \frac{\partial^2 u}{\partial t^2} &= \frac{\partial}{\partial t} \frac{\partial u}{\partial t} = -\frac{\partial}{\partial t} \left(-\alpha \frac{\partial u}{\partial x} \right) = -\alpha \frac{\partial^2 u}{\partial t \partial x} = -\alpha \frac{\partial^2 u}{\partial x \partial t} = -\alpha \frac{\partial}{\partial x} \left(\frac{\partial u}{\partial t} \right) = \\ &= -\alpha \frac{\partial}{\partial x} \left(-\alpha \frac{\partial u}{\partial x} \right) = \alpha^2 \frac{\partial^2 u}{\partial x^2} \end{aligned}$$

α is constant

the two derivatives commute

So we have:

$$u(x, t^{n+1}) = u(x, t^n) - \alpha \frac{\partial u}{\partial x}(x, t^n) \Delta t + \alpha^2 \frac{\Delta t^2}{2} \frac{\partial^2 u}{\partial x^2}(x, t^n) + \text{H.O.T.}$$

Now, we drop the H.O.T. (scaling with $O(\Delta t^3)$) and approximate

$$\frac{\partial u}{\partial x}(x_j, t^n) = \frac{1}{2h} (u_{j+1}^{n+1} - u_{j-1}^{n+1}) ; \quad \frac{\partial^2 u}{\partial x^2}(x_j, t^n) = \frac{1}{h^2} (u_{j+1}^{n+1} - 2u_j^{n+1} + u_{j-1}^{n+1})$$

The final scheme reads:

$$u_j^{n+1} = u_j^n - \frac{\alpha^2}{2} (u_{j+1}^{n+1} - u_{j-1}^{n+1}) + \frac{\alpha^2 \Delta t^2}{2} (u_{j+1}^{n+1} - 2u_j^{n+1} + u_{j-1}^{n+1})$$

Explicit Euler

Numerical Viscosity

(2) Wave Equation as 2×2 System

We noticed that the wave equation can be written as a first order system:

$$w_1 = \frac{\partial u}{\partial x} \quad w_2 = \frac{\partial u}{\partial t}$$

$$\left\{ \begin{array}{l} \frac{\partial w_1}{\partial t} = \frac{\partial^2 u}{\partial t \partial x} = \frac{\partial^2 u}{\partial x \partial t} = \frac{\partial w_2}{\partial x} \\ \frac{\partial w_2}{\partial t} = \frac{\partial^2 u}{\partial t^2} = \gamma \frac{\partial^2 u}{\partial x^2} = \frac{\partial w_1}{\partial x} \end{array} \right.$$

The two mixed derivatives commute

$$\frac{\partial w_1}{\partial t} + A \frac{\partial w_2}{\partial x} = 0$$

$$A = \begin{bmatrix} 0 & -1 \\ -\gamma^2 & 0 \end{bmatrix}$$

$$\det(\lambda I - A) = \lambda^2 - \gamma^2 = 0$$

$$\lambda_{1,2} = \pm \gamma$$

Eigenvectors:

$$A \underline{x} = \gamma \underline{x} \Rightarrow \begin{cases} -\gamma x_2 = \gamma x_1 \\ -\gamma^2 x_1 = \gamma x_2 \end{cases}$$

$$\underline{x}_1 = \begin{bmatrix} 1 \\ -\gamma \end{bmatrix}$$

$$A \underline{x} = -\gamma \underline{x} \Rightarrow \begin{cases} -\gamma x_2 = -\gamma x_1 \\ -\gamma^2 x_1 = -\gamma x_2 \end{cases}$$

$$\underline{x}_2 = \begin{bmatrix} 1 \\ \gamma \end{bmatrix}$$

$$T = \begin{bmatrix} 1 & 1 \\ -\gamma & \gamma \end{bmatrix} \quad T^{-1} = \frac{1}{2\gamma} \begin{bmatrix} \gamma & -1 \\ \gamma & 1 \end{bmatrix}$$

$$\Lambda = \begin{bmatrix} \gamma & 0 \\ 0 & -\gamma \end{bmatrix}, \quad |\Lambda| = \begin{bmatrix} \gamma & 0 \\ 0 & \gamma \end{bmatrix}$$

$$T^{-1} A T = \Lambda \Rightarrow A = T \Lambda T^{-1}$$

$$|\Lambda| = T |\Lambda| T^{-1} =$$

$$= \frac{1}{2\gamma} \begin{bmatrix} 1 & 1 \\ -\gamma & \gamma \end{bmatrix} \begin{bmatrix} \gamma & 0 \\ 0 & \gamma \end{bmatrix} \begin{bmatrix} \gamma & -1 \\ \gamma & 1 \end{bmatrix} = \frac{1}{2\gamma} \begin{bmatrix} 1 & 1 \\ -\gamma & \gamma \end{bmatrix} \begin{bmatrix} \gamma^2 & -\gamma \\ \gamma^2 & \gamma \end{bmatrix} = \begin{bmatrix} \gamma & 0 \\ 0 & \gamma \end{bmatrix}$$

Wave f.LW:

$$w_j^{n+1} = w_j^n - \frac{\gamma^2}{2} \Lambda (w_{j+1}^{n+1} + w_{j-1}^{n+1}) + \frac{\gamma^2 \Delta t^2}{2} (w_{j+1}^{n+1} - 2w_j^{n+1} + w_{j-1}^{n+1})$$

Then

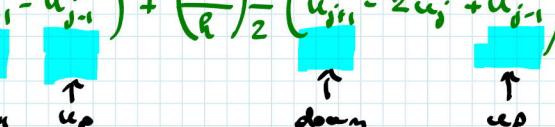
$$\frac{\partial u}{\partial t}(x_j, t^n) = w_2^n \Rightarrow w_2^n = u_j^{n-1} + 2\Delta t w_2^{n-1} \quad (\text{for example, for } n \geq 2)$$

(3) On the Boundary Conditions

We have noticed that some schemes involve points upwind and downwind.

Example

$$\text{LW: } u_j^{n+1} = u_j^n - \frac{\alpha \Delta t}{2h} (u_{j+1}^n - u_{j-1}^n) + \left(\frac{\alpha \Delta t}{h} \right)^2 \frac{1}{2} (u_{j+1}^n - 2u_j^n + u_{j-1}^n)$$

for $\alpha > 0$: 

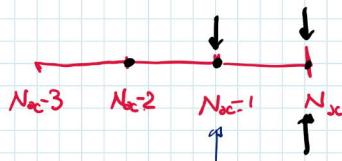
For UPW this is not happening:

$$\alpha > 0 : \quad u_j^{n+1} = u_j^n - \frac{\alpha \Delta t}{h} (u_j^n - u_{j-1}^n)$$

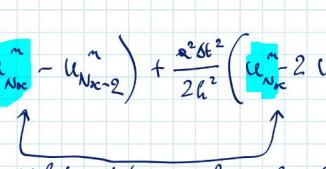
$$\alpha < 0 : \quad u_j^{n+1} = u_j^n - \frac{\alpha \Delta t}{h} (u_{j+1}^n - u_j^n)$$

This creates a problem with the B.C. because the data refer only to the inflow, we do not have data at the outflow.

For instance for $\alpha > 0$:

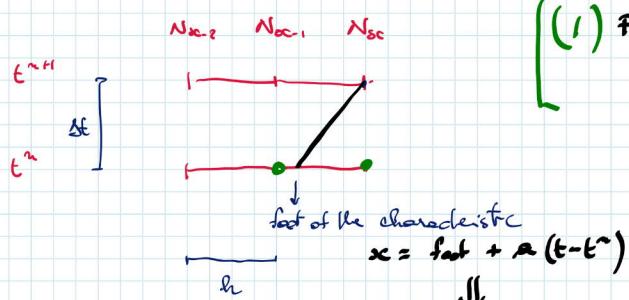


$$\text{LW: } u_{N_{xc}-1}^{n+1} = u_{N_{xc}-1}^n - \frac{\alpha \Delta t}{2h} (u_{N_{xc}}^n - u_{N_{xc}-2}^n) + \frac{\alpha^2 \Delta t^2}{2h^2} (u_{N_{xc}}^n - 2u_{N_{xc}}^n + u_{N_{xc}-2}^n)$$

 What data we have here?

There are several possible approximations.

We consider here an approach based on the analytical solution and the method of characteristics.



(1) For the equation $u_t + \alpha u_x = 0$
 $u(t^n, x_{N_{xc}}) = u(t^n, \text{foot of ch.})$

(2) For the general equation:
 $u_t + \alpha u_x + \beta u = f$
we solve:

$$\begin{cases} \frac{du}{dt} = \delta - \beta u & t \in (t^n, t^{n+1}) \\ u(\text{foot}, t^n) \text{ given} & x = \alpha(t - t^n) + \text{foot} \end{cases}$$

How do we compute $u(\text{foot}, t^n)$?

By interpolation of $[u_j^n]$.

With a piecewise linear interpolation, for instance:

$$u(\text{foot}, t^n) = u_{N_{xc}-1} \frac{x_{N_{xc}} - \text{foot}}{h} + u_{N_{xc}} \frac{\text{foot} - x_{N_{xc}}}{h}$$

We need to perform an interpolation accurate enough to maintain the accuracy.

In short, we extrapolate the missing values along the characteristics.

This approach can be extended (with some approximation) to the non-linear case.

(4) On the Boundary Conditions (2):

Reflecting conditions / Non Reflecting Conditions

When we have a system of advection equations, we found that some components travel from the left to the right, others travel back.

For instance, for the wave equation we found that one component travels forward, the other travels backward.

The compact traveling in one direction are linear combination of the physical variable \underline{u}

$$\begin{aligned}\frac{\partial \underline{u}}{\partial t} + A\underline{u} = 0 &\Rightarrow T^{-1} \frac{\partial \underline{u}}{\partial t} + T^{-1} A T T^{-1} \underline{u} = 0 \\ &\Rightarrow \frac{\partial \underline{w}}{\partial t} + D \underline{w} = 0 \\ \underline{w} &= T^{-1} \underline{u}\end{aligned}$$

Each \underline{w} :

(1) is a linear combination of the u_i ;

(2) has a well defined direction of propagation.

The physical variables \underline{u} do not have a propagation direction.

The mathematical variables \underline{w} do have a propagation direction.

If we have a 2×2 system with 1 positive & 1 negative eigenvalue, like in the wave equation:

\Rightarrow if we prescribe a condition on u_1 , or u_{12} , we are prescribing a condition on both the mathematical variables w_1 and w_{12} , so we TRIGGER COMPONENTS TRAVELING IN BOTH THE DIRECTION

In general, conditions on the physical variables are "reflecting" as they turn on components traveling in any direction.

- what if we do not want reflections ?

EXAMPLE:

$$\begin{cases} \frac{\partial u_1}{\partial t} - \frac{\partial u_2}{\partial x} = 0 \\ \frac{\partial u_2}{\partial t} - \frac{\partial u_1}{\partial x} = 0 \end{cases} \quad A = \begin{bmatrix} 0 & -1 \\ -1 & 0 \end{bmatrix} \quad \lambda^2 - 1 = 0 \quad \lambda_{1,2} = \pm 1$$

$$T = \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} \quad T^{-1} = \frac{1}{2} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}$$

$$\underline{w} = \frac{1}{2} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} \frac{u_1 - u_2}{2} \\ \frac{u_1 + u_2}{2} \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = T \underline{w} = \begin{bmatrix} w_1 + w_2 \\ -w_1 + w_2 \end{bmatrix}$$

If we prescribe a condition on either u_1 or u_2 , we affect w_1 and w_2 .

Imagine that a boundary is fictitious, so it is not a real boundary but it is just a limitation of a physical domain. For instance, assume we have a forward-propagating solution



If we prescribe conditions on \underline{u} , we introduce some components moving in both Rx directions, so we have some SPURIOUS REFLECTIONS.

To avoid artifacts induced by the boundary conditions, we need to write conditions on the \underline{w} !

In the example, this means :

Inflow: $w_1 = \text{data}$ (traveling from left to right)

Outflow: $w_2 = 0$ (no backward propagating component)

Conditions on the \underline{w} (inflow: $\frac{u_1 - u_2}{2} = \text{data}$; outflow: $\frac{u_1 + u_2}{2} = 0$) are practically more complete to prescribe, but doable.

These are called "Non-Reflecting Conditions".

Week 5 : Multi Dimensional Problems

When we go in more than 1D, we have conceptually the same problems but with many more technicalities to work out.

Before we face these topics, some notation you are probably already familiar with, but... just in case.

OPERATORS (referring to problems in 3D)

If u is a scalar function of x_1, x_2, x_3 (and possibly t):

∇u is usually referred to space variables:

$$\nabla u = \begin{bmatrix} \frac{\partial u}{\partial x_1} \\ \frac{\partial u}{\partial x_2} \\ \frac{\partial u}{\partial x_3} \end{bmatrix} \quad (\text{Gradient}) \quad \text{in 2D is a 2D vector, in 3D has 3 components}$$

If v is a vector, $\nabla \cdot v = \frac{\partial v_1}{\partial x_1} + \frac{\partial v_2}{\partial x_2} + \frac{\partial v_3}{\partial x_3}$

Clearly: $\nabla \cdot (\nabla u) = \Delta u = \frac{\partial^2 u}{\partial x_1^2} + \frac{\partial^2 u}{\partial x_2^2} + \frac{\partial^2 u}{\partial x_3^2}$

Einstein Notation: (1) an index refers to the component of a vector

EXAMPLE: $u_i = i^{\text{th}}$ component of vector u

(2) a comma + index i refers to the derivation w.r.t. x_i :

EXAMPLE: $u_{i,i} = i^{\text{th}}$ component of ∇u ($= \frac{\partial u}{\partial x_i}$)

(3) If two indices are repeated, a sum (saturation) is understood

EXAMPLE

$$u_{i,i} = \sum \frac{\partial u_i}{\partial x_i} = \nabla \cdot u \quad (= \frac{\partial u_1}{\partial x_1} + \frac{\partial u_2}{\partial x_2} + \frac{\partial u_3}{\partial x_3})$$

So, we have:

$$\Delta u = u_{i,i} \quad \begin{array}{l} \text{↑ derivation} \\ \text{index saturation} \end{array} \quad \text{outward normal unit vector}$$

Integration by parts in multiD (Green Formula):

$$\int_{\Omega} (\nabla \cdot v) u = \int_{\partial\Omega} v \cdot \underline{n} u - \int_{\Omega} v \cdot \nabla u$$

Einstein: $\int_{\Omega} v_{i,i} u = \int_{\partial\Omega} v_i n_i \cdot u - \int_{\Omega} v_i u_{,i}$

Corollary: If $v = \nabla w$ (w scalar):

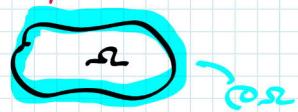
$$\int_{\Omega} \Delta w u = \int_{\partial\Omega} (\nabla w \cdot \underline{n}) u - \int_{\Omega} \nabla w \cdot \nabla u$$

Einstein: $\int_{\Omega} w_{i,i} u = \int_{\partial\Omega} w_i n_i \cdot u - \int_{\Omega} w_i u_{,i}$

NOTATION:

Space domain: Ω

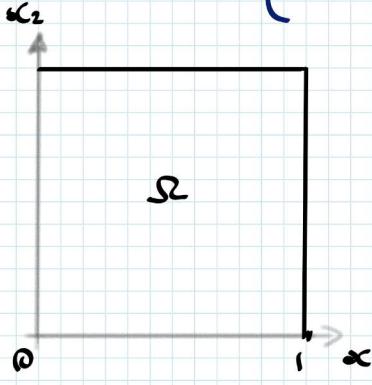
Boundary of Ω : $\partial\Omega$



(1) Poisson (Laplace) Problem in 2D

Let Ω be a square domain $((0,1) \times (0,1))$: we want to solve

$$\begin{cases} -\Delta u = f(x_1, x_2) & \text{in } \Omega \\ u(\partial\Omega) = 0 \end{cases}$$



First, recall that the problem stems from the minimization of an energy:

$$J(u) = \frac{1}{2} \int_{\Omega} |\nabla u|^2 - \int_{\Omega} fu$$

↑ elastic energy
wrk to wrk f

Physically, this may be the static deformation of a square membrane clamped at the boundaries ($\Delta u = 0$?) under the action of the force f .

If we perturb J with $u + \epsilon v$, we get: ($\epsilon \in \mathbb{R}$)

$$J(u + \epsilon v) = \frac{1}{2} \int_{\Omega} |\nabla u|^2 + \int \nabla u \cdot \nabla v + \frac{\epsilon^2}{2} \int |\nabla v|^2 - \int fv - \epsilon \int fv$$

If we take the differential:

$$\frac{1}{\epsilon} (J(u + \epsilon v) - J(u)) = \int \nabla u \cdot \nabla v + \frac{\epsilon}{2} \int |\nabla v|^2 - \int fv$$

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} (J(u + \epsilon v) - J(u)) = \int \nabla u \cdot \nabla v - \int fv$$

At the equilibrium, the differential is 0:

$$(1) \boxed{\int_{\Omega} \nabla u \cdot \nabla v = \int_{\Omega} fv \quad \forall v \text{ s.t. } v(\partial\Omega) = 0}$$

(the perturbation must be consistent with the b.c.)

With a standard application of the Green formula:

$$\begin{aligned} \int_{\Omega} \nabla u \cdot \nabla v &= \int_{\Omega} \nabla u \cdot v - \int_{\Omega} \Delta u v \\ &= 0 \text{ because } v(\partial\Omega) = 0 \end{aligned}$$

$$\Rightarrow \int_{\Omega} (-\Delta u - f)v = 0$$

v is arbitrary

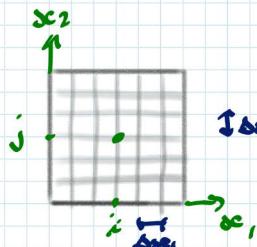
$$(2) \boxed{\text{in } \Omega: -\Delta u = f \quad (u(\partial\Omega) = 0)}$$

(1) and (2) are formally equivalent

(but in (2) u is differentiated twice, in (1) only once).

Numerical Approximation with FD

We can rely on cartesian directions :



$$u_{ij} \approx u(x_{1i}, x_{2j})$$

$$x_{1i} = i\Delta x_1, \quad x_{2j} = j\Delta x_2$$

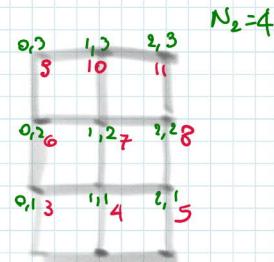
We "collocate" the problem in the vertices of the grid

$$-\Delta u(x_{1i}, x_{2j}) = f(x_{1i}, x_{2j}) = f_{ij}$$

$$-\frac{u_{i+1,j} - 2u_{i,j} + u_{i-1,j}}{\Delta x_1^2} - \frac{u_{i,j+1} - 2u_{i,j} + u_{i,j-1}}{\Delta x_2^2} = f_{ij}$$

Clearly, we need to solve a linear system, but the real issue is the organization of the solution and the state.

We need to move from a 2D space representation to a vector (1D)



Number of horizontal points : N_1

$$0 \leq i < N_1$$

$$\text{Red} = N_1 j + i \quad 0 \leq j < N_2$$

We start from 0, as this is more consistent with computer architectures.

(MATLAB starts from 1)
Python " " 0

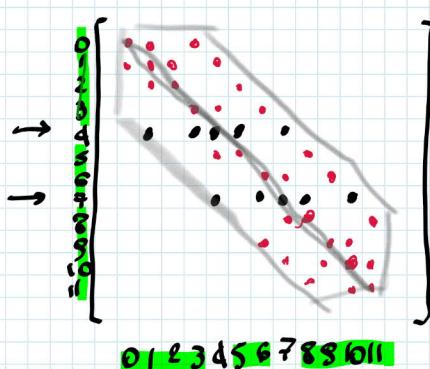
Total Number of points :

$$N_1 \times N_2 \\ (\text{from } 0 \text{ to } N_1 \times N_2 - 1)$$

We use the Red index for our unknowns :

$$u(\text{Red}(i,j)) = u_{ij}$$

Similarly, for the matrix : in



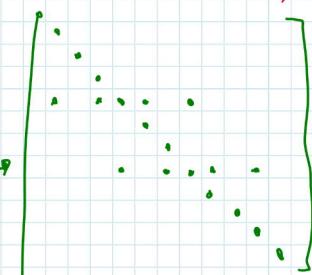
the only nodes not on the boundary are 4 and 7

However we can construct the matrix without including the B.C. The pattern of the matrix (red and black)

The matrix is banded

$$\text{Non-Diagonal} \leq 3 \quad (= N_1)$$

Pattern after B.C. :



$$\begin{aligned}
 \text{Diagonal Coefficients} \quad & \alpha_{ii} = -\frac{2}{\Delta x_1^2} - \frac{2}{\Delta x_2^2} \\
 \pm 1 \text{ coefficients:} \quad & \alpha_{i,i+1} = \frac{1}{\Delta x_1^2} \\
 \pm 3 \quad " \quad & \alpha_{i,i+3} = \frac{1}{\Delta x_2^2}
 \end{aligned}$$

After B.C.

$$\left. \begin{aligned}
 \alpha_{ii} &= 1 \\
 \alpha_{ij} &= 0 \quad j \neq i
 \end{aligned} \right\}$$

$$i \in [0, 1, 2, 3, 5, 6, 8, 9, 10, 11]$$

(1) We can construct the matrix by bands

$$\text{sparsity} \left([-3, -1, 0, 1, 3], \left[\frac{1}{\Delta x_1^2}, \frac{1}{\Delta x_2^2}, -2\left(\frac{1}{\Delta x_1^2} + \frac{1}{\Delta x_2^2}\right), \dots \right] \right)$$

(2) The nodes marked to be "boundary" are treated accordingly)

The right hand side:

$$\text{Before the B.C.: } b_i(\text{Red}(i, j)) = f(x_1(i), x_2(j))$$

$$\text{After the B.C.: } b_i[\text{Boundary Nodes}] = \text{Boundary Values}$$

Remark

$$\text{inverting the map: } j = \text{integer division} \left(\frac{\text{Red}}{N_1} \right)$$

$$i = \text{rem of integer division} \left(\frac{\text{Red}}{N_1} \right)$$

On the solution of the system

When we include the boundary conditions and look at the matrix only for the internal nodes, we have (assuming for simplicity $\Delta x_1 = \Delta x_2 = h$)

$$A = \frac{1}{h^2} \begin{bmatrix} 4 & -1 \\ -1 & 4 \end{bmatrix} \text{ which is S.P.D. } (2_1=5, 2_2=3)$$

We will see that there are features of the matrix reflecting features of the physical problem. This is critical for the solution of the linear system and, in particular, for the choice of the method.

As you can see, it's easy in multiple dimensions to have large linear systems.

Direct Methods are an option only if the matrix is not too large, otherwise the fill-in of sparse matrices can kill up the memory.

When the matrix is SPD, the Cholesky factorization is better,

because it takes half of a memory of LU

For large matrices, the only option is iterative methods.

The Book of Langtangen considers methods like Jacobi and S.O.R. These are pretty old methods.

The current choices are :

SPD Matrices : CONJUGATE GRADIENT

General Matrices : GMRES (powerful and well known)

BiCGStab (efficient but with a behavior of the error sometimes unexpected).

Remark On the Dirichlet B.C.

As you notice, the boundary nodes are now distributed everywhere in the solution, not just first and last as in the 1D case.

They need to be LABELED to be processed in the code and the sequence:

- (1) Build the system w/out Dirichlet B.C.
- (2) Process the Dirichlet B.C.

is much more efficient than the merging of the two steps.

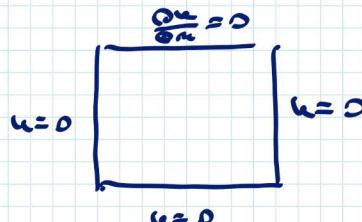
General Boundary Conditions.

The Poisson problem may feature other types of B.C.:

$$\underbrace{\mu \nabla u \cdot \underline{n}}_{\text{outward normal flux}} = \text{date} \quad (\underline{n} = \text{normal unit vector})$$

or $\mu \nabla u \cdot \underline{n} + \alpha u = \alpha u_0$ where u_0 is given.

For instance, we can prescribe on the unit square:



In this case :

$$\underline{n} = [0, 1]$$

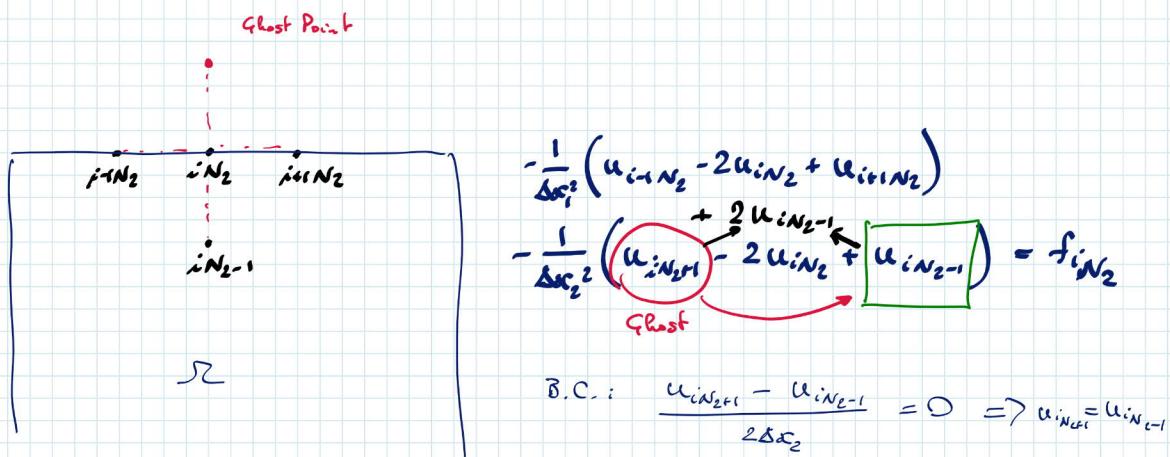
$$\text{and } \frac{\partial u}{\partial \underline{n}} \equiv \nabla u \cdot \underline{n} = \frac{\partial u}{\partial x_2}$$

To preserve the second-order nature of our scheme, we could write for the nodes with $x_2 = 1$:

$$\frac{1}{q} \left(\frac{3}{2} u_{i,N_2} - 2 u_{i,N_2-1} + u_{i,N_2-2} \right) \Leftarrow$$

$j \text{ when } x_2 = 1$

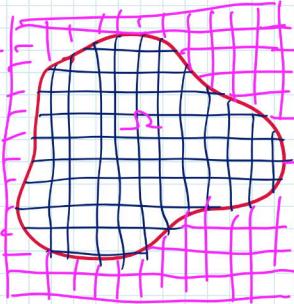
or we can write the equation also on the boundary by using ghost points:



Clearly, FD are strongly Cartesian-oriented.

Generic domains may create trouble:

Note: A possible approach is to embed the generic domain into a Cartesian frame and then solve a fictitious problem that reduces to the physical one on the domain Ω



We need many local adjustments to manage this case.

↓
We do not cover them here, as with Finite Element, they are gone!

(2) Advection-Diffusion Problems in 2D (and in 3D also)

If we want to consider a problem in multiple dimensions with diffusion and advection:

$$-\mu \Delta u + \beta \cdot \nabla u = f \quad \text{in } \Omega \quad (\mu > 0)$$

+ B.C.

where β is a vector $[\beta_1, \beta_2]$ with the direction and strength of the wind.

In detail, the problem reads:

$$-\mu \frac{\partial^2 u}{\partial x_1^2} - \mu \frac{\partial^2 u}{\partial x_2^2} + \beta_1 \frac{\partial u}{\partial x_1} + \beta_2 \frac{\partial u}{\partial x_2} = f(x_1, x_2)$$

when $\|\beta\| \gg \mu$, should we expect numerical instabilities?

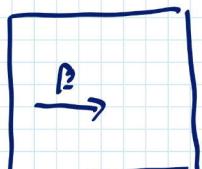
The answer is clearly positive, for the same reason of 1D: with $\|\beta\| > \mu$ the problem gets closer to change its nature, the first order term winning against the 2nd order one.

However, in 1D we noticed that upwind was equivalent to solving a numerical velocity to the problem:

$$\beta \frac{u_i - u_{i-1}}{\Delta x} = \underbrace{\beta \frac{u_{i+1} - u_{i-1}}{2\Delta x}}_{(\beta > 0)} - \underbrace{\frac{\beta \Delta x}{2} \left(\frac{u_{i+1} - 2u_i + u_{i-1}}{\Delta x^2} \right)}_{\text{numerical velocity (stabilizing)}}$$

This is not true in multiple dimensions.

To see this, let's take the case $\beta = [\beta, 0]$, so the problem reads:



$$-\mu \frac{\partial^2 u}{\partial x_1^2} - \mu \frac{\partial^2 u}{\partial x_2^2} + \beta_1 \frac{\partial u}{\partial x_1} = 0$$

Before we see numerical stabilizations, let's notice that this problem is:

- advection-diffusion along x_1 ($=$ for x_2 constant)
- purely diffusion along x_2 ($=$ for x_1 constant)

This suggests that we don't need stabilization along x_2 .

In fact, if we do UPWIND when $|\beta_1| \Delta x / 2\mu > 1$, we do:

$$-\mu \frac{u_{i+1,j} - 2u_{i,j} + u_{i-1,j}}{\Delta x_1^2} - \mu \frac{u_{i,j+1} - 2u_{i,j} + u_{i,j-1}}{\Delta x_2^2} + \beta_1 \frac{u_{i,j} - u_{i-1,j}}{\Delta x_1} = f_{i,j} \quad \text{for } \beta_1 > 0$$

$$+ \beta_1 \frac{u_{i+1,j} - u_{i,j}}{\Delta x_1} = f_{i,j} \quad \text{for } \beta_1 < 0$$

This corresponds to:

$$-\mu(1+\text{Pe}) \frac{u_{i+1,j} - 2u_{i,j} + u_{i-1,j}}{\Delta x_1^2} - \mu \frac{u_{i,j+1} - 2u_{i,j} + u_{i,j-1}}{\Delta x_2^2} + \beta_1 \frac{u_{i+1,j} - u_{i-1,j}}{2\Delta x_1} = f_{i,j}$$

We are adding viscosity only along x_1 (where really needed)



This is untouched

The "artificial viscosity" here would lead to:

$$-\mu(1+\text{Pe}) \frac{u_{i+1,j} - 2u_{i,j} + u_{i-1,j}}{\Delta x_1^2} - \mu(1+\text{Pe}) \frac{u_{i,j+1} - 2u_{i,j} + u_{i,j-1}}{\Delta x_2^2} + \beta_1 \frac{u_{i+1,j} - u_{i-1,j}}{2\Delta x_1} = f_{i,j}$$



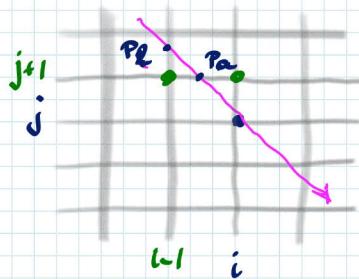
This is USELESS

In other terms, UPWIND ADDS VISCOSITY ONLY ALONG THE WIND DIRECTION (a.k.a. streamlines) while adding viscosity $\mu' = \mu(1 + \beta\epsilon)$ just adds viscosity in all the directions, including the CROSSWIND one. However, along the crosswind direction we don't need it for stabilizing, so we just add a useless numerical error.

Upwind along a generic direction

How do we implement upwind along a generic direction?

There are different strategies. We see here one, with FE we will see another one.



Let's assume to have β constant for the moment.

(streamline)

If β is constant, the upwind direction is denoted by

$$x_1 = \beta_1 z + \Delta x_{1,i} \quad \text{where } z \text{ is a}$$

$$\Delta x_2 = \beta_2 z + \Delta x_{2,j} \quad \text{parameter}$$

For $z=0$ we are in $\Delta x_{1,i}, \Delta x_{2,j}$

Assuming (just for the sake of the example) that $\beta_1 > 0, \beta_2 < 0$ we can identify the upwind point. We can find two points: one for $\Delta x_1 = \Delta x_{1,i-1}$ ($\beta_1 > 0$), the other with $\Delta x_2 = \Delta x_{2,j-1}$ ($\beta_2 < 0$)

$$P_a : \bar{x}_a = \frac{\Delta x_{2,j+1} - \Delta x_{2,j}}{\beta_2} ; \quad x_1 = \beta_1 \bar{x}_a + \Delta x_{1,i} \quad (\bar{x}_a < 0)$$

$$P_b : \bar{x}_b = \frac{\Delta x_{1,i-1} - \Delta x_{1,i}}{\beta_1} ; \quad \Delta x_2 = \beta_2 \bar{x}_b + \Delta x_{2,j}$$

Now, to make the approximation of the first derivative, we choose the point closest to $x_{1,i}, x_{2,j}$, in this case P_a .

Then:

$$\beta \cdot \nabla u \approx \beta_1 \frac{u_{1,i} - u(P_a)}{\Delta x_1} + \beta_2 \frac{u(P_a) - u_{2,j}}{\Delta x_2}$$

We have the problem that $u(P_a)$ in general is not known.

We can proceed by interpolation, for instance:

$$u(P_a) \approx u_{1,i-1,j+1} \frac{\Delta x_{1,i} - \Delta x_{1,a}}{\Delta x_1} + u_{1,i,j+1} \frac{\Delta x_{1,a} - \Delta x_{1,i-1}}{\Delta x_1}$$

using the green points for the interpolation.

As you may notice, the upwind of order 2 requires a lot of technicalities, all related to the "Cartesian" structure of the method.

With Finite Elements we will see a more general approach, still based on adding numerical viscosity, but not isotropically as done before.

Week 6 Module 2: The Finite Element Method

(1) Some Recall of Functional Analysis

Quinton: Chap. 2
Formaggia et al.: Chap. 1

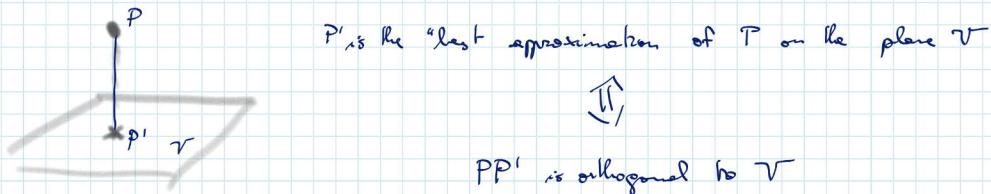
In this week we will introduce some basic tools for the analysis of differential equations of interest.

We need some concepts of functional analysis to have the right framework.

Functional Analysis, at the bottom line, is the branch of mathematics trying to extend to functions some geometrical concepts like "distance" and "orthogonality".

These concepts are extremely helpful for us:

- "distance" between two functions is related to measuring the error between the exact solution and the numerical one.
- "orthogonality" between two functions is related to the optimality of approximation:



Clearly, we do not cover many things here, Functional Analysis is a course more, but we recall here some basic tool.

Functionals: a functional is an application associating a real number to a function:

For instance, given a continuous function on an interval (a, b) ,

$$\int_a^b f(x) dx \text{ is a functional.}$$

If a functional δ is s.t., given two functions f_1, f_2 :

$$\delta(\lambda f_1 + \mu f_2) = \lambda \delta(f_1) + \mu \delta(f_2)$$

we say that the functional is linear.

Forms: suppose to have an application between a couple of functions and a real number:

$$a(u, v) \rightarrow \mathbb{R} \quad \text{for given } u \text{ and } v.$$

This is called a Form.

The form is bilinear if:

$$a(\lambda u_1 + \mu u_2, v) = \lambda a(u_1, v) + \mu a(u_2, v)$$

$$a(u, \lambda v_1 + \mu v_2) = \lambda a(u, v_1) + \mu a(u, v_2)$$

A bilinear form is symmetric when (for all the possible arguments) :

$$\alpha(u, v) = \alpha(v, u)$$

Functional Spaces

A set is a collection of objects

A "space" is a set with "operations", i.e. a set of objects that can be combined:

X is a space: $u \in X, v \in V$

if

$u \in X \quad \forall \lambda \in \mathbb{R}$

$u + v \in X \quad \forall u, v \in X$

$(\lambda u + \mu v \in X \quad \forall u, v \in X)$

$\lambda, \mu \in \mathbb{R}$

A space is "metric" if we can introduce a non-negative functional that can measure the distance between two objects:

$a \in X$
 $b \in X$

$$d(a, b) \geq 0 \text{ and}$$

$$d(a, b) = 0 \Leftrightarrow a = b$$

$$d(a, b) = d(b, a)$$

$$d(a, c) \leq d(a, b) + d(b, c)$$

With the distance, we can introduce a "topology"

A "closed" space is such that for a convergent sequence a_n , such that

$$\lim_{n \rightarrow \infty} d(a, a_n) = 0 \text{ then } a \in X.$$

A "Cauchy sequence" is a sequence such that

$$\lim_{m, n \rightarrow \infty} d(a_m, a_n) = 0$$

In general, being complete is not sufficient (but necessary) for the convergence in the space. If it is also sufficient, the space is called "complete".

A norm is a functional s.t.

$$\|f\|_X \geq 0 \text{ and } = 0 \text{ for } f = 0$$

$$\|\lambda f\|_X = |\lambda| \|f\|_X \quad \forall \lambda \in \mathbb{R} \text{ and } f \in X$$

$$\|f_1 + f_2\|_X \leq \|f_1\|_X + \|f_2\|_X$$

If a space has a norm, it is called "normed". In fact, a norm induces the following definition of distance:

$$d(a, b) = \|a - b\|_X$$

Definition: A normed complete space is called a Banach space

A scalar product $(\cdot, \cdot)_X$ is a symmetric bilinear form defined with both the arguments in a given space such that :

$$(f, f) \geq 0 \quad (= 0 \text{ only for } f=0), \quad |(f_1, f_2)| \leq (f_1, f_1)^{\frac{1}{2}} (f_2, f_2)^{\frac{1}{2}}$$

A scalar product induces a norm :

$$\|f\|_X^2 = (f, f)_X \quad (\text{e.g. } (f_1 + f_2, f_1 + f_2)_X = (f_1, f_1) + 2(f_1, f_2) + (f_2, f_2) \leq \|f_1\|^2 + 2\|f_1\|\|f_2\| + \|f_2\|^2 = ((f_1, f_1) + (f_2, f_2)))$$

So, a space with a scalar product is automatically normed.

A Hilbert space is a Banach space equipped with a scalar product.

The importance of scalar products is in the definition of ORTHOGONALITY : from a geometrical to a functional definition :

$$f \perp g \text{ in a Hilbert space} : (f, g) = 0$$

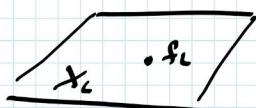
Thanks to the scalar product, we can define

- the "angle" between two functions : $\varphi(f_1, f_2) : \cos(\varphi) = \frac{(f_1, f_2)}{\|f_1\| \|f_2\|}$
- the "optimal approximation of a function in a subspace"

Let $X_L \subset X$ where X_L is "small", both Hilbert spaces :

we want to find the best approximation of f in X_L :

$\circ f$



$$(f - f_L, v_L) = 0$$

$$\forall v_L \in X_L$$

This is a consequence of the Pythagorean Theorem.

If g_L is in X_L :

$$\begin{aligned} \|f - g_L\|^2 &= (f - g_L, f - g_L) = (f - f_L + f_L - g_L, f - f_L + f_L - g_L) = \\ &\leq \|f - f_L\|^2 + 2(f - f_L, f_L - g_L) + \|f_L - g_L\|^2 = \\ &= \|f - f_L\|^2 + \|f_L - g_L\|^2 \stackrel{f_L - g_L = 0 \text{ (orthogonal)}}{\geq} \|f - f_L\|^2. \end{aligned}$$

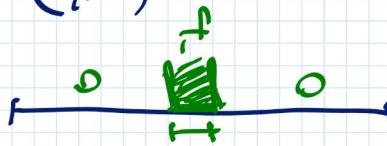
Lebesgue Integration and Some Functional Space

Differential calculus as we learned in undergraduate courses has some limitations when approaching PDE's. It's strongly related to the concept of "continuity", but we know that reality and real problems are plenty of discontinuities.

For instance, when we have a string (steady conditions):

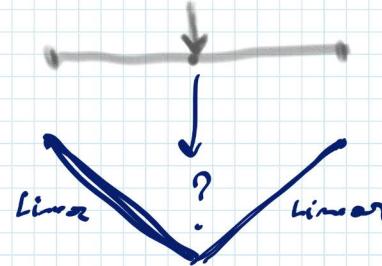
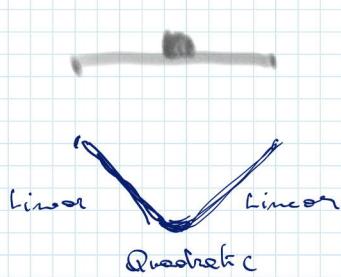
$$-u'' = f \quad \text{in } (a, b)$$

it may be discontinuous:



We can compute the solution,

but we have some formal inconsistency, particularly if we let the support of the force f (= the interval where $f \neq 0$) to have an infinitesimal length



We may have only an intuitive understanding but with some formal problems, like for the solution $u_0 + a_0 x = 0$ with a discontinuous initial condition.

$$\int [u_0(x)] \Rightarrow u(x, t) = u_0(x - at)$$

We know that this is the solution, even if formally this cannot be differentiated in space.

We need now conceptualizations for integrals and derivatives.

Lebesgue integrals are a new and somehow different way of thinking of integrals, based on a new theory of measure introduced by the French mathematician Lebesgue.

Instead of focusing on the function itself, Lebesgue integration focuses also on the " dx "

$$(L) \int_a^b f(x) dx$$

If dx is null-measure, then locally the behavior of the integral

may be "pathologic" but the function can be integrated.

The traditional and the Lebesgue integrals are largely the same (the same on regular functions) but with Lebesgue we can face functions with local lack of regularity, exactly what we need for PDE's.

Based on this, we can introduce the following spaces:
(in the sense of Lebesgue)

$$L^p(a,b) : \left\{ f : \text{(C)} \int_a^b |f|^p dx < +\infty \right\} \quad \text{with } p = 1, 2, \dots \in \mathbb{N}$$

These are Banach spaces with the norm $\|f\|_p = \left(\int_a^b |f|^p dx \right)^{\frac{1}{p}}$

(similar definition in $\Omega \subset \mathbb{R}^n$ $n \geq 1$).

$$L^\infty(a,b) : \left\{ f : \text{bounded (up to null-measure intervals) in } (a,b) \right\}$$

Banach space with $\|f\|_\infty = \max_{(a,b)} |f|$

These definitions can be extended to $(a,b) = \mathbb{R}$ or, in general, unbounded domains.

In particular: $L^2(a,b)$ is a Hilbert space with the scalar product

$$(f,g) = \int_a^b f g dx \quad (\text{Norm: } \left(\int_a^b |f|^2 dx \right)^{\frac{1}{2}}).$$

$L^2(a,b)$ or $L^2(\Omega)$ will be our "home" for most of our problems.

Notice that if we have a function $f \in L^p(a,b)$ and $g \in L^q(a,b)$ with $\frac{1}{p} + \frac{1}{q} = 1$, then

$$\text{(C)} \int_a^b f g < +\infty \quad (\text{or } fg \in L^1(a,b))$$

L^p and L^q are conjugate spaces.

Clearly, the only space conjugate with itself is L^2 , and we have

$$|(f,g)|_2 \leq \|f\|_2 \|g\|_2.$$

Distributions and distributional derivatives

The set of linear bounded functionals defined on a space X is per se a Banach space, that we call dual X' :

$$X' = \{ \delta(f) \quad \text{for } f \in X, \text{ linear and s.t. } |\delta(f)| \leq c \|f\|_X \}$$

Notice that for linear functionals bouness and continuity coincide.

Notice also that in L^2 , the dual of L^2 is isometric to L^2 .
(This is true for any Hilbert space).

This means that any linear and continuous functional in L^2 corresponds to an element in the L^2 itself (Riesz representation theorem).

The space \mathcal{D}

We define "support" of a function the interval or region of the domain of definition where the function is $\neq 0$.

In particular, if the support is a closed bounded region, we say that the function has a **COMPACT support**.

With the letter \mathcal{D} we denote the space of functions in $C^\infty(\mathbb{R}^n)$ ($n \geq 1$) with compact support.



Let introduce the dual space of \mathcal{D} , \mathcal{D}' with the following notation:

$$\mathcal{D}' \langle T, f \rangle_{\mathcal{D}}$$

$$T \in \mathcal{D}'$$

$$f \in \mathcal{D}$$

EXAMPLES:

① Let $\mathbf{1}(x)$ be the step function $= \begin{cases} 0 & \text{for } x \leq 0 \\ 1 & \text{for } x > 0 \end{cases}$

$$\langle T, f \rangle = \int_{-\infty}^{+\infty} \mathbf{1} f dx = \int_0^{+\infty} f dx$$

② The functional $\langle T_0, f \rangle = f(0)$ belongs to \mathcal{D}'

in fact $|\langle T_0, f \rangle| \leq \max_{x \in \mathbb{R}} |f(x)|$

The objects of \mathcal{D}' are called DISTRIBUTIONS or Generalized Functions.

We can introduce a differential calculus for distributions.

In fact, we formally define:

$$\langle T', f \rangle = -\langle T, f' \rangle$$

Thus it is inspired by integration by parts. In fact, if we take a L^2 function for instance, called g , we have:

$$\int_{\mathbb{R}} g' f = [gf]_{-\infty}^{+\infty} - \int_{\mathbb{R}} gf'$$

\Rightarrow because $\lim_{x \rightarrow \pm\infty} f = 0$

$$\langle g', f \rangle = -\langle g, f' \rangle$$

The definition obviously applies to higher derivatives.

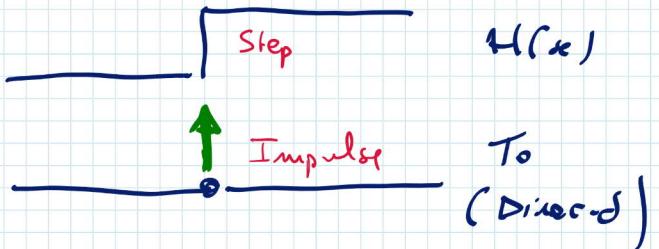
Notice that being $f \in C^\infty$, we can differentiate as many times as we want.

EXAMPLE

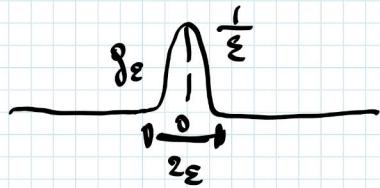
$$\int_{-\infty}^{+\infty} H' f = - \int_{-\infty}^{+\infty} H f' = - \int_0^{+\infty} f = -[f]_0^{+\infty} = \delta(0) = \langle T_0, f \rangle$$

T_0 is called Dirac- δ and it is the mathematical description of an impulse.

It is not a function in the traditional sense



Dirac- δ is the limit of functions:



$$\lim_{\epsilon \rightarrow 0} g_\epsilon = \delta$$

$$\text{where } \int_{\mathbb{R}} g_\epsilon = 1 \quad (\forall \epsilon)$$

The definition of derivatives in the sense of distributions is very general and with a wide range of applications.

$$L^2 \subset \mathcal{D}'$$

FROM NOW ON, WHEN WE DIFFERENTIATE A FUNCTION IN L^2 , WE INTEND A DISTRIBUTIONAL DERIVATIVE.

Sobolev Spaces

Once we have defined the new concept of derivative
that doesn't require a strong constraint like continuity,
we can characterize functions that can be differentiated

Space $H^k(a, b)$ ($\text{or } H^k(\Omega)$) is the space of functions f s.t. the belong to L^2
together with all the derivatives up to the order k .

$$H^k(a, b) = \left\{ f \text{ s.t. } f, f', \dots, f^{(k)} \in L^2 \right\}$$

In multiple dimensions, the derivatives are all the derivatives of order k .

$$\frac{\partial^{\alpha_1 + \dots + \alpha_n} f}{\partial x_1^{\alpha_1} \partial x_2^{\alpha_2} \dots \partial x_n^{\alpha_n}} \in L^2 \quad (*)$$

with $\sum \alpha_i = k$

These spaces are called SOBOLEV SPACES (Russian mathematicians) and are a special case of the more general definition:

$$W^{k,p} = \left\{ f \text{ s.t. } f, f', \dots, f^{(k)} \in L^p \right\}$$

$(H^k = W^{k,2})$.

There are many properties of spaces H^k . Let's just recall some of them:

$\Rightarrow H^k(a, b)$ is a Hilbert space with the scalar product

$$(f, g)_{H^k} = \sum_{i=0}^k (f^{(i)}, g^{(i)})_{L^2}$$

differentiation index

(in 1D, in multi-dimension
the sum is extended to all the
derivatives in the form (*)).

Specifically, $H^1(a, b)$ has the norm:

$$\|f\|_{H^1}^2 = \|f\|_{L^2}^2 + \|f'\|_{L^2}^2$$

Clearly $H^1 \subset L^2$ so $\|f\|_{H^1}^2 \geq \|f\|_{L^2}^2$ (as obvious from the definition).

In general $H^{k+1} \subset H^k \subset \dots \subset L^2$.

Notice that the definition of Sobolev spaces H^k can be extended to n real and negative.

The Space H'_0

In general, the value of a function in a point or on a line (in 2D) or a surface (in 3D) doesn't matter because L^2 is defined up to null-measure sets and a point/line/surface has null measure in 1D/2D/3D.

However, in PDE's we prescribe boundary conditions, so we need to understand in what sense we can prescribe BC.

This has been clarified by the introduction of the concept of TRACE. TRACE is an operator that maps a function in H' to its value at the boundary and, in general, reduces regularity:

$$\gamma: H'(\Omega) \rightarrow H^{1/2}(\partial\Omega)$$

is the (continuous) trace operator.

From now on, BC will be intended in the sense of traces.

The space $H'_0(\Omega)$ is the space of H' functions with null trace on $\partial\Omega$.

This will be very important for our problems.

The scalar product and norm of H'_0 are like H' . However we have a fundamental property in H'_0 :

Poincaré' INEQUALITY: $\exists C > 0$ s.t. $\forall f \in H'_0(\Omega)$:

$$\|f\|_{L^2} \leq C \|\nabla f\|_{L^2}$$

This allows to write the following sequence of inequalities:

$$\|\nabla f\|_{L^2}^2 \leq \|f\|_{L^2}^2 + \|\nabla f\|_{L^2}^2 \leq (1+C^2) \|\nabla f\|_{L^2}^2$$

\Rightarrow In $H'_0(\Omega)$, the norm of ∇f in L^2 and of f in H' are equivalent

This is true also if the function has only a portion of the boundary (Γ with measure > 0) where it vanishes.

The reason why the Poincaré inequality has to be extended in H'_0 is that in general we may have a function with the norm \geq of the gradient: the constants! However, in H'_0 the only constant possible is 0.

The space of linear and continuous functionals on H'_0 is denoted by H^{-1} and we have that

$$H^{-1} \text{ contains } L^2 \text{ contains } H'_0$$

Space-Time Functional Spaces:

For time-dependent problems we have spaces that describe the two dependencies.
For instance

$L^p(0,T; H^k(\Omega))$ is the space of functions such that:

the norm $\|f\|_{H^k}(\cdot) \in L^p(0,T)$ [in short: $L^p(H^k)$].

We will use $L^2(H^0)$ and $L^\infty(L^2)$

Embedding Theorems

There is a relation between belonging to H^k and L^p with $p \neq 2$.
Also, there is a relation between H^k and the traditional continuity.

These relations are summarized by a sequence of theorems called "embedding theorems".

If $u \in H^s$ in $\Omega \subset \mathbb{R}^n$

$$0 < 2s < n : \quad H^s \subset L^q \quad 1 \leq q \leq q^* = \frac{2n}{n-2s}$$

(continuous)
embedding

$$2s = n : \quad H^s \subset L^q \quad \forall q \quad 1 \leq q < +\infty$$

$$2s > n : \quad H^s \subset C^0(\bar{\Omega})$$

Also

$$H^s \subset C^m(\bar{\Omega}) \quad \text{if} \quad s > m + \frac{n}{2}$$

Notice the role of the number of space-dimensions!

This explains why some results of well-posedness (e.g. in fluid mechanics) are available only in 2D and not in 3D.

Question

Let u be a function in $H^1(\Omega)$ s.t. $u(\partial\Omega) = 1$

Can I define the space

$$H_1(\Omega) = \{ f \in H^1, f(\partial\Omega) = 1 \} \quad ?$$

NO! It's not a space! In fact: $f_1 \in H^1, f_2 \in H^1 \Rightarrow f_1 + f_2 \notin H^1$.

(2) The Poisson Problem in 1D

We have already noted that the classical problem

$$-u'' = f \quad x \in (0,1)$$

$$u(0) = u(1) = 0$$

is in fact the result of the minimization of

$$J = \frac{1}{2} \int_0^1 (u')^2 - \int_0^1 f u$$

On the way, we found that this correspond to solving the problem:

$$\int_0^1 u' v' = \int_0^1 f v \quad \forall v \text{ s.t. } v(0) = v(1) = 0.$$

Let's give a more rigorous formulation.

With the notation already introduced we have:

find $u \in H_0^1(0,1)$: $\int_0^1 u' v' dx = \int_0^1 f v dx \quad \forall v \in H_0^1(0,1)$

where $f(x) \in L^2(0,1)$.

This is what we call the "weak" formulation of the problem.

Notice that:

\Rightarrow the BC are encoded in the space, H_0^1

$\Rightarrow f \in L^2$ is not the most general space, as a matter of fact we can take $f \in H^{-1}(0,1)$

\Rightarrow each term is well defined:

$$\int_0^1 u' v' dx \quad u' v' \text{ is summable}$$

\uparrow \uparrow
 $u' \in L^2$ $v' \in L^2$

$$\int_0^1 f v$$

\uparrow
 $v \in H_0^1$

$\Rightarrow f$ can be a continuous-linear function on H_0^1 .

\Rightarrow We can prescribe $f = \delta(\frac{x}{2})$ for the problem depicted on the right

Is this problem well-posed?

let's give a general formulation

$$\boxed{\text{Find } u \in V : \alpha(u, v) = f(v) \quad \forall v \in V} \quad (1)$$

- where :
- (i) $\alpha(\cdot, \cdot)$ is a bilinear form: $V \times V \rightarrow \mathbb{R}$
 - (ii) $f(\cdot)$ is a linear-continuous function on V : $V \rightarrow \mathbb{R}$.

The LAX-MILGRAM LEMMA :

For the problem (1), let's assume:

- (A) $\alpha(\cdot, \cdot)$ continuous: $\exists M: |\alpha(u, v)| \leq M \|u\|_V \|v\|_V$
- (B) $f(\cdot)$ " : $\exists C: |f(v)| \leq C \|v\|_V$

and in particular:

$$(C) \exists \gamma > 0: \alpha(u, u) \geq \gamma \|u\|_V^2 \quad \forall u \in V$$

COERCIVITY OF THE BILINEAR FORM

(A), (B), (C) are sufficient to conclude that:

- (I) u exists
- (II) u is unique
- (III) u depends continuously on the date

Conclusion (I) can be found in functional analysis books.

Conclusion (II) by contradiction:

If we have u_1, u_2 solving the same problem:

$$H_0' \ni u_1, u_2 : \alpha(u_1 - u_2, v) = 0 \quad \forall v \in H_0'$$

Then: for $v = u_1 - u_2$:

$$\alpha \|u_1 - u_2\|_V^2 \leq \alpha(u_1 - u_2, u_1 - u_2) \leq 0$$

$$\downarrow \\ \|u_1 - u_2\|_V = 0 \Rightarrow u_1 - u_2 = 0 \quad (\text{contradiction})$$

Also: We have for $v = u$:

$$\alpha \|u\|_V^2 \leq \alpha(u, u) \leq C \|u\|_V \quad \Rightarrow \quad \|u\|_V \leq \frac{C}{\alpha}$$

(stability bound for the solution)

$$(iii) \text{ Let } u_p : \quad \mathcal{Q}(u_p, v) = \delta(v) + \underbrace{\delta\mathcal{J}(v)}_{\downarrow} \quad \xrightarrow{\text{Continuous perturbation}} \quad \forall v \in H^1_0(0, 1)$$

For $v = u_p - u$:

$$\alpha \|u_p - u\|_v^2 \leq \mathcal{Q}(u_p - u, u_p - u) \leq \delta C \|u_p - u\|_v^2$$

$$\Rightarrow \|u_p - u\|_v \leq \frac{\delta C}{\alpha}$$

so, for $\delta \rightarrow 0$, $\|u_p - u\|_v \rightarrow 0$.

For our problem, Lax-Milgram Lemma can be applied \Rightarrow EXERCISE

WARNING: The Lax-Milgram lemma is a **SUFFICIENT** (but not NECESSARY) condition. If it is not applicable, then the problem can be still well-posed.

LAX-MILGRAM Disease = to think that LM is also necessary 😊

The LM Lemma is a particular case of a more general result called BNB (Babuška, Nečas, Banach) Theorem (that, conversely, is necessary & sufficient.)

More General Boundary Conditions

(1) Non-homogeneous Dirichlet conditions

$$u(0) = d_0, \quad u(1) = d_1$$

As we noted, we cannot just introduce a set of functions fulfilling the B.C., because this set is not a space.

The simple solution is to introduce a LIFTING or EXTENSION. Let $l(x)$ to be an arbitrary function of H^1 s.t.

$$l(0) = d_0, \quad l(1) = d_1$$

Then, let $\tilde{u} = u - l$, so that:

$$\tilde{u}(0) = \tilde{u}(1) = 0$$

$$\text{and } \mathcal{Q}(\tilde{u}, v) = \underbrace{\mathcal{Q}(u, v) - \mathcal{Q}(l, v)}_{\text{known}}$$

So, formally we can write the problem for \tilde{u} : find $\tilde{u} \in H^1_0$ s.t.

$$\mathcal{Q}(\tilde{u}, v) = \mathcal{Q}(u, v) - \mathcal{Q}(l, v) \quad \forall v \in H^1_0$$

and then we reconstruct the solution $u = \tilde{u} + l$.

Since the solution is unique, we know that regardless of the specific lifting

the solution will be the same.

How do we select the lifting function $l(x)$?

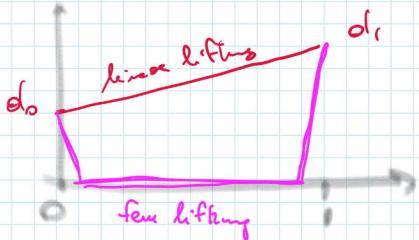
In 1D there is a natural choice, linear interpolation:

$$l(x) = d_0 + (d_1 - d_0)x \quad x \in (0, 1)$$

By the way, this is such that $l''=0$, so:

$$-u'' = -\ddot{u} - l'' = -\ddot{u} = f$$

For the numerical approximation, there are better choices.



We will see that using the piecewise polynomial lifting in the picture (magenta) the finite element method can manage the lifting automatically, with no additional effort.

(2) Neumann (or Natural) conditions

Let's assume that in $x=1$ we have the boundary condition $u'(1)=r \in \mathbb{R}$

In the "strong" form we have:

$$-u'' = f \quad x \in (0, 1), \quad u(0) = 0, \quad u'(1) = r$$

For the weak form, instead of using the minimization argument, we use a formal practical approach.

(a) multiply the equation by v : $-u''v = fv$

(b) integrate over the domain $(0, 1)$: $-\int_0^1 u'' v = \int_0^1 fv$

(c) integrate by parts:

$$\begin{aligned} & -[u'v]_0^1 + \int_0^1 u'v' = \int_0^1 fv \\ & -u'(1)v(1) + u'(0)v(0) \end{aligned}$$

(d) manage the boundary conditions:

- take $v(0) = 0$ (like u vanishes in 0)

- notice that $u'(1) = r$

Weak form: let $H'_D = \{f \in H' \text{ s.t. } f(0) = 0\}$:

$$\text{find } u \in H'_D \text{ s.t. } \int_0^1 u'v' = \int_0^1 fv + rv(1) \quad \forall v \in H'_D$$

Also in this case, the Lax-Milgram lemma can be applied and we have well poshers.

Notice that HOMOGENEOUS Neumann conditions ($r=0$) lead to:

$$\int_0^1 u' v' = \int_0^1 f v \quad H_D'$$

so these conditions do not require any modification of the functional and of the functional space either.

They are called "do-nothing" conditions.

REMARK

Let's consider this problem:

$$-u'' = f \quad u'(0) = r_0 \quad u'(1) = r,$$

The weak formulation reads:

$$\text{find } u \in H^1(0,1) \text{ s.t.}$$

$$\int_0^1 u' v' = \int_0^1 f v + r_1 v(1) - r_0 v(0)$$

Is this problem well posed?

- Direct answer: u occurs in the equation and the B.C. only under differentiation, so it can be replaced by $u + \text{constant}$ and the equation is still the same \Rightarrow NO UNIQUENESS.
- Functional analysis answer: if you test the coercivity of the bilinear form:

$$\alpha(u, v) = \int_0^1 u' v'$$

Notice that:

$$\alpha(u, u) = \int_0^1 u'^2 = \|u'\|_{L^2}^2$$

The Poincaré inequality states that $\|u'\|_{L^2}$ is equivalent to $\|u\|_{H^1}$ in H^1_0 and H_D^1 , not in H^1 .

\Rightarrow The Lax-Milgram lemma doesn't apply.

(3) Robin or Mixed Conditions:

Let's assume the condition:

$$u'(1) + \sigma u(1) = r \quad (\text{and } u(0) = 0)$$

The weak formulation reads:

find $u \in H_0^1$ s.t.

$$\begin{aligned}\int_0^1 u' v' &= \int_0^1 f v + [u' v]_0^1 = \\ &= \int_0^1 f v + (\nu(u(1)) - \nu(u(0))) v(1) \quad \forall v \in H_0^1\end{aligned}$$

or:

$$\int_0^1 u' v' + \nu(u(1)) v(1) = \int_0^1 f v \quad \forall v \in H_0^1.$$

(3) The Poisson Problem in 2D (or 2+D)

In this case, we noticed that the problem reads:

$$-\Delta u = f \quad \text{in } \Omega \quad u(\partial\Omega) = 0$$

The weak formulation with our formal approach follows the same guidelines:

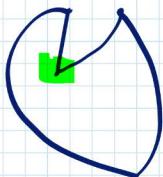
$$(1) \quad -\Delta u v = f v \quad \text{with } v(\partial\Omega) = 0$$

$$(2) \quad -\int_{\Omega} \Delta u v = \int_{\Omega} f v$$

$$(3) \quad -\int_{\Omega} \Delta u v = - \int_{\Omega} (\nabla u \cdot \nabla v) + \int_{\Omega} \nabla u \cdot \nabla v$$

$$(4) \Rightarrow \text{Find } u \in H_0^1(\Omega) \text{ s.t. } \int_{\Omega} \nabla u \cdot \nabla v = \int_{\Omega} f v \quad \forall v \in H_0^1(\Omega)$$

Now, the Lax-Milgram Lemma can be applied only with some constraints on Ω :

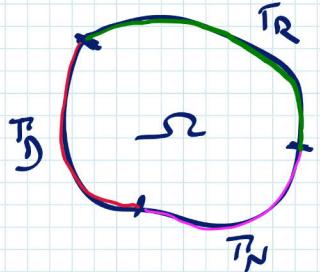


Concavities can undermine the continuity of the bilinear forms and the functionals.

In general, for polygonal domains, the well posedness is guaranteed.

We can consider more general sets of boundary conditions.

Let's consider a general problem with different boundary conditions.



$$-\Delta u = f \quad \text{in } \Omega$$

$$u(\partial\Omega) = g_d \quad \text{on } \Gamma_D$$

$$\nabla u \cdot \underline{n} = g_N \quad \text{on } \Gamma_N$$

$$\nabla u \cdot \underline{n} + \sigma u = g_R \quad \text{on } \Gamma_R$$

$$\partial\Omega = \overline{\Gamma_D \cup \Gamma_N \cup \Gamma_R}$$

$$\Gamma_N \cap \Gamma_D = \emptyset$$

$$\Gamma_N \cap \Gamma_R = \emptyset$$

$$\Gamma_R \cap \Gamma_D = \emptyset$$

(1st) Let's introduce a lifting, i.e. a regular function l in Ω , s.t.

$$l(\partial\Omega) = g_d.$$

Then, defining $\tilde{u} = u - l$:

$$-\Delta(\tilde{u} + l) = f \Rightarrow -\Delta \tilde{u} = f \quad \text{in } \Omega$$

$$\tilde{u}(\Gamma_D) = 0$$

$$\nabla \tilde{u} \cdot \underline{n} = \underline{d}_N - \nabla l \cdot \underline{n} \quad \text{on } \Gamma_N$$

$$\nabla \tilde{u} \cdot \underline{n} + \sigma \tilde{u} = \underline{d}_R - \nabla l \cdot \underline{n} - \sigma l$$

Then, with the usual procedure, notice that the Green formula here applies as:

$$-\int_{\Omega} \Delta u v = - \int_{\Gamma_D} \nabla u \cdot \underline{n} v - \int_{\Gamma_N} \nabla u \cdot \underline{n} v - \int_{\Gamma_R} \nabla u \cdot \underline{n} v$$

So, we have: Find $\tilde{u} \in H^1_{\Gamma_D}(\Omega)$ ($= f \in H^1$ s.t. $l(\Gamma_D) = 0$) s.t.

$$\int_{\Omega} \nabla \tilde{u} \cdot \nabla v + \int_{\Gamma_R} \sigma \tilde{u} v = \int_{\Omega} f v + \int_{\Omega} \Delta l v + \int_{\Gamma_N} (\underline{d}_N - \nabla l \cdot \underline{n}) v + \int_{\Gamma_R} (\underline{d}_R - \nabla l \cdot \underline{n} - \sigma l) v$$

$$\text{Alternative: } \int_{\Omega} f v - \int_{\Omega} \nabla l \cdot \nabla v + \int_{\Gamma_N} \underline{d}_N v + \int_{\Gamma_R} (\underline{d}_R - \sigma l) v \quad \forall v \in H^1_{\Gamma_D}(\Omega)$$

(if we decide to apply the Green formula to l too).

The analysis of this problem where we have several boundary conditions at the same time is not trivial \Rightarrow An entire book by P. Grisvard is dedicated to this.

REMARK 1: You may notice that the same problem has two different weak formulations depending on the treatment of the lifting.

Keep in mind that, in general, after numerical approximation they may lead to slightly different numerical results.

REMARK 2: More in general, it is important to stress that a weak formulation does not refer to an equation, but to the entire Boundary Value problem: DIFFERENT BOUNDARY CONDITIONS LEAD TO DIFFERENT WEAK

FORMULATIONS.

REMARK 3: In all the formulations above you can check that if $f \in L^q(\Omega)$ (or also to the dual of $H_{\Gamma_D}^1$), then the terms of the weak formulation are all well defined.

(A) More General Problems.

In general, let us consider the following problem:

$$-\nabla \cdot (\mu(x_1, x_2) \nabla u) + \beta(x_1, x_2) \cdot \nabla u + \sigma(x_1, x_2) u = f \quad \text{in } \Omega$$

- with
- $u|_{\Gamma_D} = 0$ (for non-homogeneous state we need a lifting)
 - $\mu(x_1, x_2) \nabla u = d_N \quad \text{on } \Gamma_N$
 - $\mu(x_1, x_2) \nabla u + \gamma(x_1, x_2) u = d_R \quad \text{on } \Gamma_R$

Let's construct the weak formulation formally. Then, we will discuss the well-posedness.

Green Formula

$$-\int_{\Omega} \nabla \cdot (\mu \nabla u) v = -\int_{\partial\Omega} \mu \nabla u \cdot \nu v + \int_{\Omega} \mu \nabla u \cdot \nabla v$$

Find $u \in H_{\Gamma_D}^1(\Omega)$ s.t.

$$\int_{\Omega} \mu \nabla u \cdot \nabla v + \int_{\Omega} (\beta \cdot \nabla u) v + \int_{\Omega} \sigma u v + \int_{\Gamma_R} \gamma u v = \int_{\Omega} f v + \int_{\Gamma_N} d_N v + \int_{\Gamma_R} d_R v$$

$$\forall v \in H_{\Gamma_D}^1$$

To analyze the well-posedness, we need first to guarantee that each term makes sense ($=$ is $<+\infty$).

Let's assume $f \in L^2(\Omega)$

$$\int_{\Omega} \mu \nabla u \cdot \nabla v$$

\downarrow \downarrow
 $\in L^2$ $\in L^2$
EL'

Since $(\nabla u \cdot \nabla v) \in L^1$
 Then we need
 $\mu \in L^\infty(\Omega)$

When we have

$$\int_{\Omega} \varphi \psi$$

with $\varphi \in L^p$, to have this term bounded we need $\psi \in L^{p'}$
 with $p' = \left(1 - \frac{1}{p}\right)^{-1} = \frac{p}{p-1}$
 If $\varphi \in L^1$, $\psi \in L^\infty$.

$$\int_{\Omega} \beta \cdot \nabla v \quad \left[\begin{array}{l} \text{For the embedding theorem, } v \in L^p \text{ where } p \text{ depends} \\ \text{on the number of dimensions.} \end{array} \right]$$

↓

In 2D and 3D $v \in L^6(\Omega)$, so we can assume
 $\beta \in (L^3)^m$ (each component of the vector
belongs to L^3).

$$\int_{\Omega} \sigma u v \quad \left[\begin{array}{l} \text{We need } \sigma \in L^{\frac{3}{2}}(\Omega) \text{ at least.} \\ \text{Diagram: } \begin{array}{c} \downarrow \quad \downarrow \\ H^1 \quad H^1 \\ \underbrace{\quad}_{\in L^3} \end{array} \end{array} \right]$$

$$\int_{\Gamma_R} \gamma u v \quad \rightarrow \text{we need } \gamma \in L^2(\Gamma_R) \\ \left[\begin{array}{l} \downarrow \quad \downarrow \\ H^{\frac{1}{2}} \quad H^{\frac{1}{2}} \\ (\text{traces}) \\ \underbrace{\quad}_{L^1} \end{array} \right]$$

So, as you can see, the "collection" in the most general spaces is not easy.

In general, in the applications, it's common to have:

$$u \in L^\infty(\Omega), \quad \beta \in (L^\infty(\Omega))^m, \quad \sigma \in L^\infty(\Omega), \quad \gamma \in L^\infty(\Gamma_R)$$

So we are good to go, because $L^\infty(\Omega)$ is in all the spaces indicated above.

Lax-Milgram - based well posedness

Assuming that all the functions $\mu, \beta, \sigma, \gamma$ and t are regular enough, so the bilinear forms and the functionals are all continuous, UNDER WHAT CONDITIONS, the bilinear form:

$$a(u, v) = \int_{\Omega} \mu \nabla u \nabla v + \int_{\Omega} (\beta \cdot \nabla u) v + \int_{\Omega} \sigma u v + \int_{\Gamma_R} \gamma u v$$

is COERCIVE?

If we limit for simplicity to the case $\mu > 0, \sigma \geq 0, \gamma \geq 0$ it is immediately realized that the 1st, 3rd, 4th of the bilinear form

$a(u, u)$ are positive.

We need to focus on:

$$\int_{\Omega} (\beta \cdot \nabla u) u \quad \left(\text{Einstein notation: } \int_{\Omega} \beta_i u_i u_j u_j \right)$$

Notice that we can write:

$$\frac{1}{2} \int_{\Omega} \beta \cdot \nabla (u^2) \underset{\text{Green}}{=} \frac{1}{2} \int_{\Omega} \beta \cdot \nabla u^2 - \frac{1}{2} \int_{\Omega} (\nabla \cdot \beta) u^2$$

This term for general boundary conditions is complicated.

For Dirichlet conditions ($u(\partial\Omega) = 0$), the first term is 0 and we conclude that

$$\text{if } (\nabla \cdot \beta) \leq 0 \Rightarrow \int_{\Omega} (\beta \cdot \nabla u) u \geq 0.$$

As you may notice, the Lax-Milgram analysis gets complicated!!!

REMARK

In the construction of the weak formulation, we can integrate by parts also the term:

$$\begin{aligned} \int_{\Omega} \beta \cdot \nabla u v &= \int_{\Omega} \beta \cdot \nabla u v - \int_{\Omega} \nabla \cdot (\beta v) u = \\ &= \int_{\Omega} \beta \cdot \nabla u v - \int_{\Omega} (\nabla \cdot \beta) v u - \int_{\Omega} \beta \cdot \nabla v u \\ &\quad \uparrow \\ &\quad v|_{\Gamma_D} = 0 \text{ so this reduces to } \Gamma_N \cup \Gamma_R \end{aligned}$$

Since the term on the boundary changes, this approach makes sense for a different boundary condition, for instance:

$$\mu \nabla u \cdot n + \beta \cdot \nabla u = \tilde{d}_N \quad \text{on } \Gamma_N$$

or

$$\mu \nabla u \cdot n + \beta \cdot \nabla u + g u = \tilde{d}_R \quad \text{on } \Gamma_R$$

} in some applications, these are, in fact, the physical conditions to prescribe.

Recently the Remark a couple of pages back, FOR THE SAME EQUATION WE MAY HAVE DIFFERENT WEAK FORMULATIONS SUITABLE FOR DIFFERENT BOUNDARY CONDITIONS.

Week 7

Questions: Chap. 3
QV - Chap. 5

The Galerkin Method

In this chapter, we have a general view and introduce the Galerkin method for the generic problem:

$$\text{find } u \in V : \quad a(u, v) = f(v) \quad \forall v \in V$$

where $a(\cdot, \cdot)$ is a bilinear form: $V \times V \rightarrow \mathbb{R}$

$f(\cdot)$ a linear and continuous functional: $V \rightarrow \mathbb{R}$

We assume in general that V is an infinite dimensional space.

In the FD method, one approximation was based on the idea of approximating the differential operators (with incremental quotients), after "collocating" the problem in special points.

With the Galerkin method we change completely our perspective.

The approximation is not in the differential operators, but in the search space.

Let V_h be a FINITE DIMENSIONAL SUBSPACE of V , where each function v_h can be represented as the linear combination of BASIS FUNCTIONS:

$$v_h = \sum_{j=1}^{N_h} c_j \varphi_j$$

where N_h is the dimension of V_h .

Then, we construct the APPROXIMATE PROBLEM

$$\text{find } u_h \in V_h \text{ s.t. } a(u_h, v_h) = f(v_h) \quad \forall v_h \in V_h$$

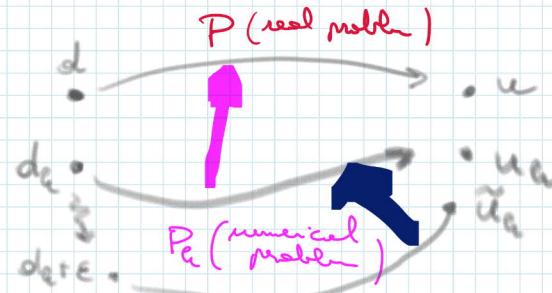
This is the Galerkin approximation of our problem.

Before we look into the specific construction of V_h , we can state some general properties of this approach.

You recall that for an approximation of a differential problem we introduced the concepts of

- **CONSISTENCY** $P \xrightarrow{\epsilon \rightarrow 0} P$

- **STABILITY** $\tilde{u}_h \xrightarrow{\epsilon \rightarrow 0} u_h$



For the consistency, if we construct $\nabla_h \subset \nabla$ is such a way that $\nabla_h \xrightarrow{h \rightarrow 0} \nabla$, we are off.

But we can characterize the consistency error in a better way:

In general, if we pretend the exact solution to fulfill the numerical problem, we have a RESIDUAL:

$$\begin{aligned} P(u; d) &= 0 \\ P_h(u_h; d_h) &= 0 \quad (\text{in exact arithmetic}) \end{aligned} \quad \left. \begin{array}{l} P_h(u, d_h) \neq 0 \\ \downarrow \\ \text{Truncation error} \end{array} \right\}$$

You may recall that **CONSISTENCY** means that the truncation error vanishes with h .

EXAMPLE:

$$\frac{d(e^x)}{dx} = e^x$$

$$\frac{e^{x+h} - e^x}{h} \neq e^x \quad \text{but } e^{x+h} = e^x + e^x h + \frac{e^x}{2} h^2 + \dots$$

$$\text{error: } \frac{e^x}{2} h \xrightarrow{h \rightarrow 0} 0$$

The Galerkin method is more than consistent, it is **STRONGLY CONSISTENT**: the truncation error is 0.

P: $u \in \nabla : \alpha(u, v) = f(v) \quad \forall v \in \nabla$

P_h : $u_h \in \nabla_h \subset \nabla : \alpha(u_h, v_h) = f(v_h) \quad \forall v_h \in \nabla_h$

If we test P_h with u , we have (for $\nabla_h \subset \nabla$):

$$\alpha(u, v_h) = f(v_h)$$

because $v_h \in \nabla$

so we have **STRONG CONSISTENCY** and:

$$\alpha(u - u_h, v_h) = 0$$

$$\forall v_h \in \nabla_h$$

This is a fundamental relation.

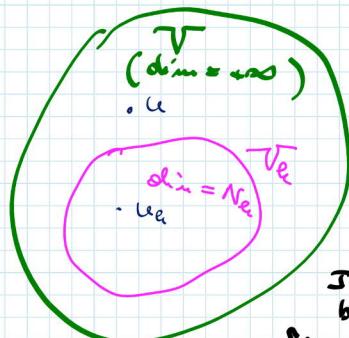
Let's recall a definition:

A bilinear form is **SYMMETRIC** if:

$$\alpha(u, v) = \alpha(v, u) \quad \forall u, v \in \nabla$$

Now, let's assume that our problem P (find $u \in \nabla : \alpha(u, v) = f(v)$ for all $v \in \nabla$) falls under the LAX-MILGRAM lemma.

Under this assumption, it is immediate to prove that also P_h is well posed! *



If the bilinear form is already in ∇ , it is in ∇_h as well

Now, we have

$$|\varphi(u, v)| \leq M \|u\|_V \|v\|_V$$

$$\varphi(u, u) \geq \alpha \|u\|_V^2 \Rightarrow \text{positive definite} \quad \varphi(u, u) \leq M \|u\|_V^2$$

In addition, let's assume $\varphi(u, v)$ to be Symmetric.

Then, we can state that $\varphi(\cdot, \cdot)$ defines a scalar product having all the features to be so.

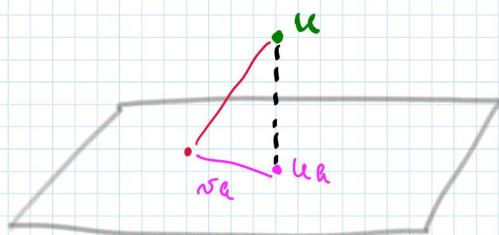
So, we can define a norm too,

$$\|u\|_E = (\varphi(u, u))^{\frac{1}{2}}$$

and this norm is equivalent to the norm in V , because

$$\alpha \|u\|_V^2 \leq \varphi(u, u) \leq M \|u\|_V^2$$

The relation $\varphi(u - u_h, v_h) = 0$ is therefore an orthogonality relation: this states that for the topology induced by $\varphi(\cdot, \cdot)$, u_h is the "best" approximation of u in V_h :



$$\varphi(u - u_h, v_h) = 0$$

$$u - u_h \perp V_h$$

- The strong consistency is independent of the symmetry.
- The geometrical interpretation (highly suggestive) holds only for symmetric forms

However, the strong consistency leads to an important result.

$$\begin{aligned} \alpha \|u - u_h\|_V^2 &\leq \varphi(u - u_h, u - u_h) = \varphi(u - u_h, u - w_h + w_h - u_h) = \\ &\quad \uparrow \quad \uparrow \quad \downarrow \\ &\quad V \quad V \quad \forall w_h \in V_h \quad \in V_h \\ &= \varphi(u - u_h, u - w_h) + \varphi(u - u_h, w_h - u_h) \\ &\quad \cancel{=} 0 \\ &\leq M \|u - u_h\| \|u - w_h\| \end{aligned}$$

for strong consistency



$$\|u - u_h\|_V \leq \frac{M}{\alpha} \|u - w_h\|_V$$

Since v_h is arbitrary, we can take w_h to be the "best" approximation of u in V_h in the topology of V :

$$w_h : \|u - w_h\|_V \leq \|u - v_h\|_V \quad \text{if } v_h \in V_h$$

or

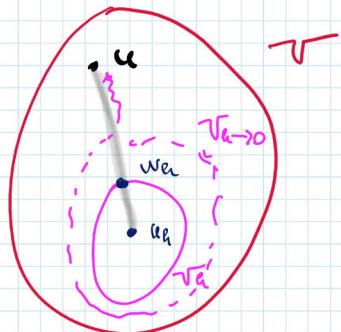
$$w_h : \inf_{v_h \in V_h} \|u - v_h\|_V = \|u - w_h\|_V$$

So we have that :

$$\|u - u_h\|_V \leq \frac{1}{\alpha} \inf_{v_h \in V_h} \|u - v_h\|_V$$

This is called Céa Lemma.

The Céa Lemma states that the Galerkin solution is not necessarily the best one, but converges to the exact one with the same rate of the best approximation.



If $a(\cdot, \cdot)$ is symmetric, u_h is the "best" in the norm $\|\cdot\|_a$ and it is possible to prove that this norm is equivalent to the norm in V .

$$u_h : \|u - u_h\|_a = \inf_{v_h \in V_h} \|u - v_h\|_a$$

$$\text{In general : } \|u - u_h\|_V \leq c \inf_{v_h \in V_h} \|u - v_h\|_V$$

(2) Our Possible Choices for the space V_h .

We can now become more specific and specify some possible choices for V_h .

In general, we can construct the solution as a polynomial function.

If we do this, the theory of polynomial approximation (MATH 516) can help understanding the term

$$\inf_{v_h \in V_h} \|u - v_h\|_V$$

One possible option is using Gaussian Interpolation Theory to construct the functions u_h . As well known from Interpolation Theory, the collocation of the interpolation nodes is critical for the convergence of the interpolation. Gaussian interpolation is an option leading to a class of methods called SPECTRAL METHODS.

The book of Langtangen - Mardal "Introduction to Numerical Methods for

various problems treat this topic in Chapter 2 and 4.

Here, we follow a different approach (we will explain why later on).

We select V_h in this way. Let's start with (1) problems.

The Finite Element Method in 1D.

(1) We take the interval ($[0, 1]$ or $[a, b]$ in general) and split it into subintervals



Notice that the distance $|x_{j+1} - x_j|$ does not need to be constant.
We set

$$h = \max_j |x_{j+1} - x_j| \quad T_h \text{ denotes the set of nodes } \{x_j\}$$

For simplicity, we will just assume $|x_{j+1} - x_j| = h \quad \forall j$
(uniform reticulation)

(2) We construct V_h as the space of ^{continuous} piecewise polynomial functions with degree p on the reticulation:

$$V_h(T_h) \equiv \left\{ v_h \in C^0([0, 1]) \text{ s.t. } v_h(I_j) \in P^p \right\}$$

(3) We represent each function v_h with a Lagrange piecewise polynomial basis:

$$v_h = \sum_{j=1}^{N_h} c_j \varphi_j(x)$$

where

$$\varphi_j(x) \in V_h \text{ and } \varphi_j(x_i) = \delta_{ij}$$

We will see examples very soon.

(4) How do we use all of this?

We want to solve

$$\text{find } u \in V_h : a(u_h, v_h) = f(v_h) \quad \forall v_h \in V_h$$

(4.1) Do we need to test against all the functions

$$v_h \in V_h ?$$

No! Clearly, if a generic function is:

$$u_h = \sum v_j \varphi_j(x) \quad (\star)$$

it is enough to use only the functions φ_i :

$$\text{Find } u_h : \alpha(u_h, \varphi_i) = f(\varphi_i) \quad \forall i = 1, \dots, N_h \quad (\square)$$

In fact, all the other equations are a linear combination of these, thanks to the linearity of the arguments in u_h and the (x).

(4.2) But what is (\square) from a practical point of view?

Also the solution u_h reads:

$$u_h = \sum v_j \varphi_j(x)$$

so that we obtain:

$$\text{Find } \underline{v} = [v_j]_{j=1 \dots N_h} \text{ s.t. } \sum_{j=1}^{N_h} v_j \alpha(\varphi_j, \varphi_i) = f(\varphi_i) \quad \forall i = 1 \dots N_h$$

Now, set

$$A = \left[\alpha_{ij} = \alpha(\varphi_j, \varphi_i) \right] \quad N_h \times N_h \text{ matrix (called STIFFNESS MATRIX)}$$

$$\underline{v} = \begin{bmatrix} v_1 \\ \vdots \\ v_{N_h} \end{bmatrix} \quad \underline{f} = \begin{bmatrix} f(\varphi_1) \\ \vdots \\ f(\varphi_{N_h}) \end{bmatrix} \quad \text{(called LOAD VECTOR)}$$

Then, it is easily verified that we have a linear system:

$$A \underline{v} = \underline{f}$$

Problem (\square) is just a linear system: we know how to solve it!

What are the features of this linear system?

① First of all, we already proved that the problem P_0 is well posed, so we argue that the matrix is non-singular.

Also notice that the coercivity implies:

$$\alpha(u_h, u_h) \geq \alpha \|u_h\|^2$$

Also, by contradiction, assume that $\exists \lambda = 0$ eigenvalue of A ($\Rightarrow A$ is singular).

Then let u_0 be the associated eigenvector and $u_{h0} = \sum u_{j0} \varphi_j$.

Then

$$\alpha \|u_{h0}\|^2 \stackrel{\text{(coercivity)}}{\leq} \alpha(u_{h0}, u_{h0}) \stackrel{(\lambda=0 \text{ eigenvalue})}{=} u_0^T A u_0 = u_0^T 0 = 0$$

CONTRADICTION: all $\lambda \neq 0 \Rightarrow A$ is not singular.

② Notice that if $\alpha(\cdot, \cdot)$ is symmetric, it follows that

A is s.p.d.

③ On the sparsity pattern of A we'll discuss in the description of specific finite elements.

(1) \Rightarrow the linear system is non-singular

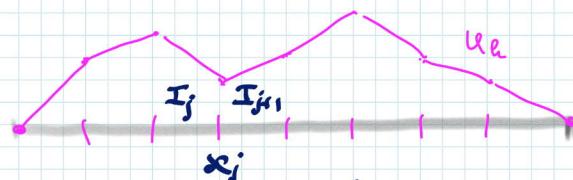
(2) \Rightarrow CG or Cholesky can be used

The question to answer is the CONVERGENCE RATE.

However, before this, let's see some particular examples.

Linear Finite Elements

In this case, we approximate the solution with a continuous piecewise linear function. We call this case as P1.



The Lagrange basis functions take the form:

$$\varphi_j(x) = \begin{cases} 0 & x \notin I_j \cup I_{j+1} = [x_{j-1}, x_{j+1}] \\ \frac{x - x_{j-1}}{x_j - x_{j-1}} & x \in I_j = [x_{j-1}, x_j] \\ \frac{x_{j+1} - x}{x_{j+1} - x_j} & x \in I_{j+1} = [x_{j+1}, x_j] \end{cases}$$

(if the mesh is uniform
 $x_j - x_{j-1} = x_{j+1} - x_j = h$)

Notice that :

- (1) φ_j are nonzero only on two elements (or intervals)
We say that φ_j has a small support
- (2) the coefficients u_j in $u_h = \sum u_j \varphi_j(x)$ have
the meaning of :

$$u_j = u_h(x_j)$$

This is the consequence of having chosen a Lagrange basis function.

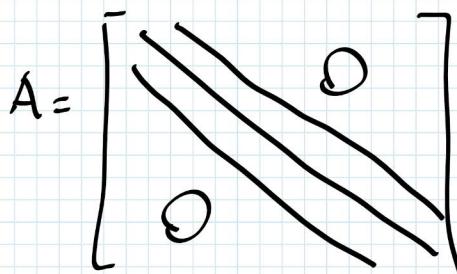
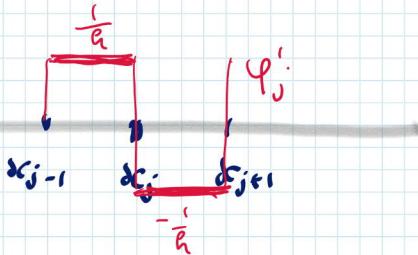
Using Lagrange polynomials ($\varphi_j(x_i) = \delta_{ij}$ (Kronecker delta)) is not the only possible choice. We will see this later.

Property (1) has an important consequence :

Let's look at the entries of A :

$a(\varphi_j, \varphi_i)$ in the Laplace operator is

$\int (\varphi_j' \varphi_i')$ \Rightarrow if $j \neq i, i+1, i-1$ this integral is automatically 0, because the support of φ_j' and the one of φ_i' do not intersect.



$$a_{ij} = 0 \text{ if } j \neq i, i+1, i-1$$



A is tridiagonal
(like for FD)

A note on the number of degrees of freedom (DOF):

In P1 we have a coefficient for each vertex:

If we have N intervals, we have $N+1$ nodes ($N_h = N+1$)
but if we have 2 Dirichlet conditions, the real size of the system reduces to $N_h - 2$ ($= N-1$) equations (in theory ... practice is different).

If we adopt an element-wise perspective,

{ On each element we want a linear function
 { A linear function is specified by 2 coefficients

Locally, we have 2 DOF needed to allow the computation of the unique function.

Choosing "Lagrange" interpolation means that we use the 2 DOF to prescribe the values of the function in two nodes (if we knew the function u , we would say that we "interpolate" u in the 2 nodes)

Choosing the vertexes guarantees the continuity of the functions

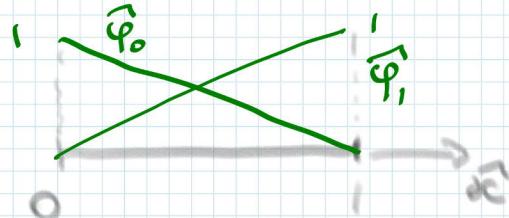
↳ This is important because this guarantees that $T_h \in P1 \subset H^1$ ("conforming" finite element).

The Reference Element

With the element-wise perspective, we may introduce an auxiliary tool.

Let's assume for a while that we have only 1 element (and do not link to the B.C. for now).

We write this in the reference variable \hat{x} :



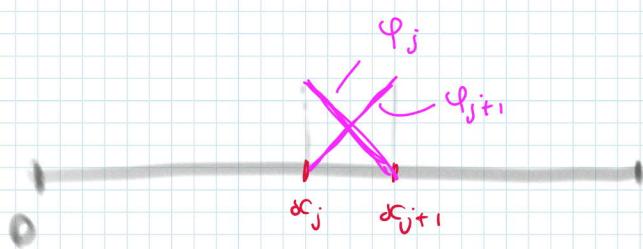
Here we have two basis functions

$$\hat{\varphi}_0 = 1 - \hat{x}$$

$$\hat{\varphi}_1 = \hat{x}$$

$$\text{On this element, } \Delta_{\text{loc}} = \begin{bmatrix} \varphi_0(\hat{\varphi}_0, \hat{\varphi}_0) & \varphi_0(\hat{\varphi}_1, \hat{\varphi}_0) \\ \varphi_0(\hat{\varphi}_0, \hat{\varphi}_1) & \varphi_1(\hat{\varphi}_1, \hat{\varphi}_1) \end{bmatrix}$$

Notice that we have the same situation on each element of a generic reticulation:



$$\varphi_j(x) = \hat{\varphi}_0(\hat{x})$$

$$\text{where } x = x_j + \hat{x}(x_{j+1} - x_j)$$

$$(\hat{x} = \frac{x - x_j}{x_{j+1} - x_j})$$

$$\varphi_{j+1}(x) = \hat{\varphi}_1(\hat{x})$$

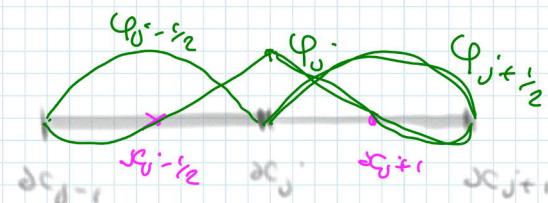
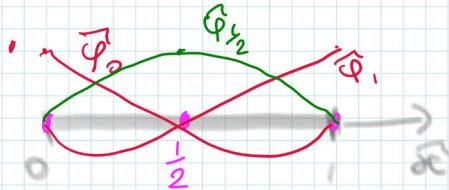
While in 1D this use of the Reference Finite Element is redundant, it is critical in multidimensions.

Quadratic Finite Elements

Let's adopt the element-wise perspective.

On each interval, we want to approximate the solution with a quadratic function:

$$T_h = P_2 = \{ v_h \in C^0([0,1]) \text{ s.t. } v_h(x_j) \in P_2 \}$$



We need 3 DOF on each interval to be unisolvant (= uniquely find the function).

We use three Lagrangian DOF, so that the coefficient of the expansion is the value of the solution in the point

We choose as interpolation points:

- the vertices (so we have continuity between two intervals)
- the mid-point (for symmetry and to have better numerical properties)

The corresponding Lagrange polynomials are depicted in the figure:

The total number of DOF (3 per interval but the ones in the vertices coincide on pairs of intervals)

$$N_h = 3 \times N - (N-1) = 2N+1$$

\uparrow
vertices

or : $\left. \begin{array}{l} 1 \text{ dof per vertex} \\ 1 \text{ dof per interval} \end{array} \right\} N_h = \overbrace{N+1}^{\text{vertices}} + N = 2N+1$

If we include the ^{Doubt}B.C. : $N_h = 2N-1$

$$\widehat{\varphi}_0 = \frac{(\widehat{x} - \frac{1}{2})(\widehat{x} - 1)}{\frac{1}{2}} = 2(\widehat{x} - \frac{1}{2})(\widehat{x} - 1) = 2(\frac{1}{2} - \widehat{x})(1 - \widehat{x})$$

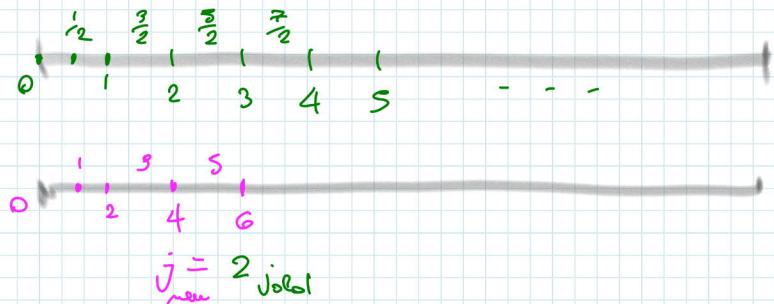
$$\widehat{\varphi}_{1/2} = \frac{\widehat{x}(\widehat{x} - 1)}{-\frac{1}{4}} = 4\widehat{x}(1 - \widehat{x})$$

$$\widehat{\varphi}_1 = \frac{\widehat{x}(\widehat{x} - \frac{1}{2})}{-\frac{1}{2}} = 2\widehat{x}(\frac{1}{2} - \widehat{x})$$

Then we can φ_j as the map with $x = x_j + \widehat{x}(x_{j+1} - x_j)$

The function $\widehat{\varphi}_{1/2}$ (and $\varphi_{0 \pm 1/2}$) are called "bubbles," because they vanish at the boundary.

Notice that in this case, if we re-number the degrees of freedom:



we have that :

$\left\{ \begin{array}{l} \text{the support of } \varphi_{2k} \text{ is } I_{2k} \cup I_{2k+1} \\ \text{the support of } \varphi_{2k+1} \text{ is } I_{2k+1} \end{array} \right.$

Consequently :

$$Q(\varphi_j, \varphi_i) = \begin{cases} \text{i even : } 0 & \text{for } j \neq i, i \pm 1, i \pm 2 \\ \text{i odd : } 0 & \text{for } j \neq i, i \pm 1 \end{cases}$$

In this case, the matrix is penta-diagonal, alternating 5 and 3 non-zero entries

$$0 \begin{bmatrix} \vdots & & & & & \\ & \ddots & & & & \\ & & \ddots & & & \\ & & & \ddots & & \\ & & & & \ddots & \\ & & & & & \ddots \end{bmatrix} \quad \leftarrow \text{then we can prescribe the B.C. here}$$

\leftarrow and here

Different Choices for the Basis

For what we have done so far, if we compute the solution with P1 and then we want to recompute the solution with P2, we have to recompute the matrix.

$$x_0 \quad x_1 \quad x_2 \quad \dots \quad x_{N-1} \quad x_N \quad x_{N+1} \quad \dots \quad x_j$$

Assuming to have two Dirichlet conditions, we have and $N = \# \text{ of intervals}$

$$P1 : N_d = N - 1 \text{ dofs}$$

$$P2 : N_d = 2N - 1 \text{ dofs}$$

Size of matrices:

$$A_{P1} : (N-1) \times (N-1) \text{ with } (N-1) + 2(N-2) = 3N-5 \text{ nonzero entries}$$

$$A_{P2} : (2N-1) \times (2N-1) \text{ with } 3N + 5(N-1) = 8N-5 \text{ nonzero entries}$$

TASK: construct a basis "prone to recycling" so that A_{P1} is a submatrix of A_{P2} .

In the reference element, we can take:

$$\left. \begin{array}{l} \hat{\varphi}_0 = 1 - \hat{x} \\ \hat{\varphi}_1 = \hat{x} \end{array} \right\} \text{ same as } P1 \quad \hat{\varphi}_{\frac{1}{2}} = \hat{\varphi}_0 \hat{\varphi}_1$$

This basis is unisolvent, any quadratic function can be written as:

$$p_2(\hat{x}) = c_0 \hat{\varphi}_0 + c_1 \hat{\varphi}_1 + c_{\frac{1}{2}} \hat{\varphi}_{\frac{1}{2}}$$

$$\text{where } c_0 = p_2(0), \quad c_1 = p_2(1)$$

while $c_{\frac{1}{2}}$ has no precise physical meaning.

$$\text{for instance: } p_2 = \hat{x}^2 + \hat{x} + 1 =$$

$$= 1 - \hat{x} + 3\hat{x} - (1 - \hat{x})\hat{x}$$

$$\text{or in general: } p_2 = \alpha \hat{x}^2 + \beta \hat{x} + \gamma =$$

$$= \gamma(1 - \hat{x}) + (\alpha + \beta + \gamma)\hat{x} - \alpha(1 - \hat{x})\hat{x}$$

When we construct the matrix, we may write:

$$A = \begin{bmatrix} \alpha(\hat{\varphi}_0, \hat{\varphi}_0) & \alpha(\hat{\varphi}_0, \hat{\varphi}_1) & \alpha(\hat{\varphi}_{\frac{1}{2}}, \hat{\varphi}_0) \\ \alpha(\hat{\varphi}_0, \hat{\varphi}_1) & \alpha(\hat{\varphi}_1, \hat{\varphi}_1) & \alpha(\hat{\varphi}_{\frac{1}{2}}, \hat{\varphi}_1) \\ \alpha(\hat{\varphi}_0, \hat{\varphi}_{\frac{1}{2}}) & \alpha(\hat{\varphi}_1, \hat{\varphi}_{\frac{1}{2}}) & \alpha(\hat{\varphi}_{\frac{1}{2}}, \hat{\varphi}_{\frac{1}{2}}) \end{bmatrix}$$

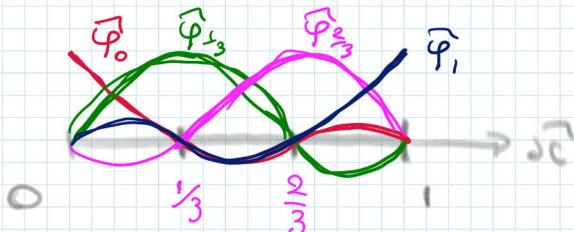
The submatrix highlighted in yellow is the P1 matrix!!!

We have a "nested" structure that may be helpful in case of an automatic selection of the degree (p -adaptivity).

This hierarchical basis is just an example of non-Lagrangian basis and it is available in FEM1D (Advanced Mode).

Cubic Finite Elements

We can proceed rapidly now



$$\hat{\varphi}_0 = \frac{(\hat{x} - \frac{1}{3})(\hat{x} - \frac{2}{3})(\hat{x} - 1)}{(-\frac{1}{3})(-\frac{2}{3})(-1)}$$

$$\hat{\varphi}_1 = \frac{\hat{x}(\hat{x} - \frac{1}{3})(\hat{x} - \frac{2}{3})}{(1 - \frac{1}{3})(1 - \frac{2}{3})}$$

$$\hat{\varphi}_{\frac{1}{3}} = \frac{\hat{x}(\hat{x} - \frac{2}{3})(\hat{x} - 1)}{\frac{1}{3}(\frac{1}{3} - \frac{2}{3})(\frac{1}{3} - 1)}$$

$$\hat{\varphi}_{\frac{2}{3}} = \frac{\hat{x}(\hat{x} - \frac{1}{3})(\hat{x} - 1)}{\frac{2}{3}(\frac{2}{3} - \frac{1}{3})(\frac{2}{3} - 1)}$$

Lagrangian DOF:

value of the function in:

$$0, \frac{1}{3}, \frac{2}{3}, 1$$

$$T_C = P_3 \equiv \left\{ v_C \in C^0([0, 1]) : v_C(I_j) \in P^3 \quad \forall I_j \right\}$$

Number of DOFs ($N = \# \text{intervals}$)

$$\underbrace{4 \times N}_{\text{of 4 dof per interval}} - (N-1) [-2 \text{ for b.c.}] = 3N+1 \quad [-2 \text{ for b.c.}]$$

continuity at
the vertices

$$\text{or } N+1 \text{ vertices} + 2 \times N \quad [-2 \text{ for b.c.}] \\ (2 \text{ int of 4 dof per interval})$$

$$= 3N+1 \quad [-2 \text{ for b.c.}]$$

You can verify that with the proper numbering, the pattern is a sequence with 7, 4, 4 non zero entries. (7 for rows with an index multiple of 3 and 4 for the others).

Convergence Rate Theory

The starting point is the Céa Lemma :

$$\|u - u_h\|_V \leq C \inf_{w \in V_h} \|u - w\|_V.$$

Let's use the interpolation theory for piecewise polynomial interpolation to quantify the right hand side.

Theorem

Let's assume $f \in H^{s+1}([a, b])$, $s \geq 1$ and let $\{x_i\}$ a partition of nodes such that $x_i + \alpha_j = f(x_{i+j})$ and set $h = \max |x_{i+1} - x_i|$.

Then, if f_h is the piecewise polynomial interpolation of order p (i.e., f_h :

$$f_h(x_i) = f(x_i) \quad \text{in the d.o.f.}$$

(vertices for $p=1$)

(vertices and element dof for $p \geq 1$)

$$f_h(I_i) \in P^p(I_i)$$

then, $\exists C > 0$ s.t.

$$\|f - f_h\|_{H^r(a,b)} \leq C h^r \left\| \frac{\partial^{(s)} f}{\partial x^{(s)}} \right\|_{L^2(a,b)}$$

$$\text{with } r = \min(p, s)$$

Also

$$\|f - f_h\|_{L^2(a,b)} \leq C h^{r+1} \left\| \frac{\partial^{(s)} f}{\partial x^{(s)}} \right\|_{L^2(a,b)}.$$

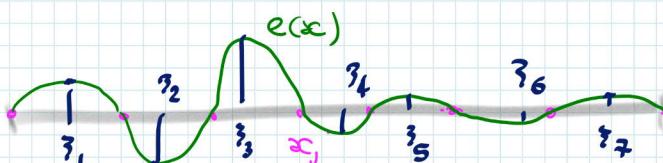
PROOF

We consider only the case $f \in H^2$, $p=1$ ($\text{so } s=1$).

$f_h(x_i) = f(x_i)$ means that $e(x) \equiv f(x) - f_h(x)$ has $N+1$ zeros (where N is the number of intervals).

The Rolle's theorem implies that in each interval there exists at least 1 point ξ_j s.t.

$$e'(\xi_j) = 0$$

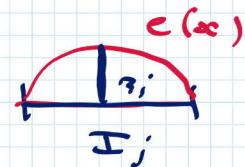


Also, notice that in each interval $f''_k = 0$ so:

$$e''(x) = f''(x) - 0 = f''(x)$$

Then:

$$e'(x) = \int_{z_j}^x e''(\sigma) d\sigma \quad \text{for } x \in I_j. \quad (\text{in fact, } e'(z_j) = 0 \text{ with this formula})$$



$$|e'(x)| \leq \left(\int_{z_j}^x 1 \right) \left(\int_{z_j}^x (e'')^2 \right)^{\frac{1}{2}} \leq h^{\frac{1}{2}} \|e''\|_{L^2(I_j)}$$

$$|e'(x)|^2 \leq h \|e''\|_{L^2(I_j)}^2 \quad x \in I_j.$$

$$\|e'\|_{L^2(I_j)}^2 = \int_{I_j} |e'(x)|^2 dx \leq h \|e''\|_{L^2}^2 \int_{I_j} 1 dx = h^2 \|e''\|_{L^2(I_j)}^2$$

$$\Rightarrow \|e'\|_{L^2(a,b)}^2 = \sum_j \|e'\|_{L^2(I_j)}^2 \leq h^2 \sum_j \|e''\|_{L^2(I_j)}^2 = h^2 \|e''\|_{L^2(a,b)}^2$$

i.e.,

$$\boxed{\|e'\|_{L^2(a,b)} \leq h \|e''\|_{L^2(a,b)} = h \|f''\|_{L^2(a,b)}} \quad (\#)$$

Now, we have:

$$|e(x)| = \left| \int_{x_j}^x e'(x) dx \right| \leq \left| \int_{x_j}^x \left(\int_{x_j}^x (e')^2 \right)^{\frac{1}{2}} dx \right| \leq h^{\frac{1}{2}} \|e'\|_{L^2(I_j)}$$

$e(x_j) = 0$

$$\|e\|_{L^2(I_j)}^2 = \int_{I_j} e^2 dx \leq h \|e'\|_{L^2(I_j)}^2 \int_{I_j} 1 dx = h^2 \|e'\|_{L^2(I_j)}^2$$

$$\Rightarrow \|e\|_{L^2(a,b)}^2 = \sum_j \|e\|_{L^2(I_j)}^2 \leq h^2 \sum_j \|e'\|_{L^2(I_j)}^2 = h^2 \|e'\|_{L^2(a,b)}^2 \leq C h^4 \|f''\|_{L^2(a,b)}^2$$

i.e.

$$\boxed{\|e\|_{L^2(a,b)} \leq C h^2 \|f''\|_{L^2(a,b)}} \quad (\#)$$

$$(I + II) \Rightarrow \|e\|_{H^1}^2 = \|e\|_{L^2}^2 + \|e'\|_{L^2}^2 \leq C(h^4 + h^2) \|f''\|_{L^2}^2$$



$$\|e\|_{H^1} \leq \tilde{C} h \|f''\|_{L^2}$$

Notice that $\|f''\|_{L^2(\Omega, \mathbb{R})}$ is also called the semi-norm H^2
(like a norm without the property $|u|_{H^2} = 0 \Leftrightarrow u = 0$)

$$\|e\|_{H^1} \leq \tilde{C} h \|f\|_{H^2}.$$

		solution regularity (H^s)						
		0	1	2	0	1	2	
FE degree	0	Conv	1	1	1	Conv	2	2
	1	Conv	1	2	2	Conv	2	3
	2	Conv	1	2	3	Conv	2	3
	3	Conv	1	2	3	Conv	2	3

↓

H^1 L^2

This table provides the guide to the selection of the most appropriate FE degree as a function of the regularity.

In fact, from the Céa Lemma:

$$\|u - u_h\|_{H^1} \leq C \inf_{w_h \in V_h} \|u - w_h\|_{H^1} \leq C \|u - u_h^*\|_{H^1}$$

where u_h^* is the interpolant of u on the dof.

From here, we notice that selecting high-order FE pays only for regular solutions, otherwise the convergence rate is the same of low-order FE.
So, the boxes outlined in green are somehow "optimal" choices.

A last piece of information is given by the following Theorem.

Theorem

The condition number of the stiffness matrix scales with h^{-2} , i.e.:
 $\text{cond}(A) \propto h^{-2}$

This means that when we refine the mesh $h \rightarrow \frac{h}{2}$ we expect $\text{cond}(A)$ to be multiplied by a factor of 4, regardless the degree of FE.

A FINAL REMARK

If we want to increase the quality of a FE solution, we have two options.

h -strategy: REFINEMENT THE MESH

p -strategy: INCREASE THE DEGREE [Hierarchical Basis very useful here]

The second strategy works only for regular solutions.

The first one increases the computational costs as stated by this last theorem.

The two strategies can be

- automated (ADAPTIVITY: we need an a-posteriori error estimator)

- combined (hp -adaptivity)

If you are interested in this topic, see Quadrilateral Chap. 4

Talking about ADAPTIVITY, M. Fornasier (great mathematician) stated that a mesh should not be static but an "unknown" of a FE solution.

Week 8

Finite Elements in 2+D for Elliptic Problems

We want to extend our procedure to the 2+D case.
The backbone of our approach is the same.

$$\text{Find } u : \quad \alpha(u, v) = f(v) \quad \forall v \in V$$

We introduce $V_h \subset V$ with $\dim(V_h) = N_h < +\infty$
and assume that

$$u_h = \sum_{j=1}^{N_h} c_j \varphi_j(\mathbf{x})$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

Then, the Galerkin formulation reads:

$$\boxed{\text{Find } u_h \in V_h \quad \alpha(u_h, v_h) = f(v_h) \quad \forall v_h \in V_h}$$

and this is equivalent to

$$\text{Find } \mathbf{c} = [c_j] : \quad \sum_{j=1}^{N_h} \alpha(\varphi_j, \varphi_i) c_j = f(\varphi_i) \quad i = 1 \dots N_h$$

In its algebraic form this problem reads:

$$\mathbf{A} \mathbf{c} = \mathbf{b}$$

$\left[\alpha(\varphi_i, \varphi_i) \right]$ (Stiffness matrix) $\left[\alpha(\varphi_i) \right]$ (Load vector)

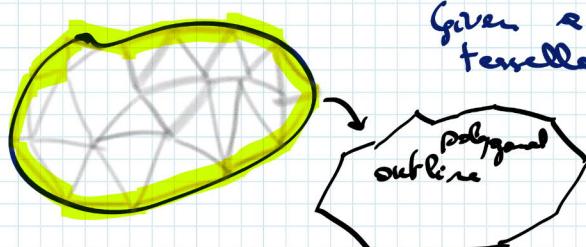
Δ inherits the features of $\alpha(\cdot, \cdot)$ as in the 1D case.

Finite Elements

The basic idea of FE is to construct the space $V_h \subset V (= H_0^1)$ as a **piecewise polynomial** approximation.

One important difference with the 1D case is the construction of the "pieces" or of the MESH

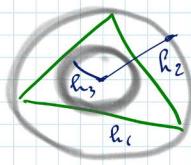
Triangular Meshes in 2D



Given a generic domain, we can think of a tessellation with triangles.

Triangles are really versatile and can cover regions with complex shapes. However, we have some **WARNINGS**.

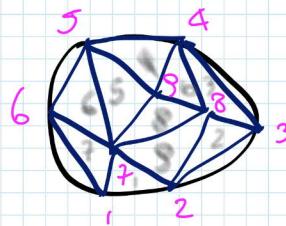
(1) Denote by T_h the tessellation with Triangles $\{T_k\}$ $k=1\dots N$
 h denotes a characteristic length of the triangles.
A possible choice is the longest edge, or the radius of
the inscribed circle or the circumscribed



(2) A mesh, in general, consists of

VERTEXES	{
EDGES (2 vertices)	
TRIANGLES (3 vertices, 3 edges)	

Now, the way a mesh file is organized differs for different softwares. However, you find this information.



List of Vertices : $x_1, x_2 \dots$ (2 in the angle)

Pattern of Elements :

1	2	7
2	3	8
3	4	8
4	5	9
5	7	9
6	5	6
7	1	7
8	7	8
9	2	8

The order is
always
COUNTER CLOCK
WISE

List of Boundary Edges

1	2
2	3
3	4
4	5
5	6
6	1

} + Labels

Additional information may be added to identify the neighbor elements :

Element 1) 9 7 (-1 = boundary : - boundary edge)

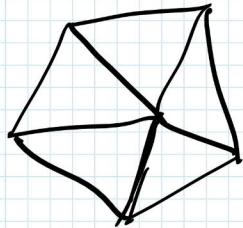
Element 2) 3 9 (-2)

(3) If we want continuous functions, we have some restrictions.
Why we want continuous functions?

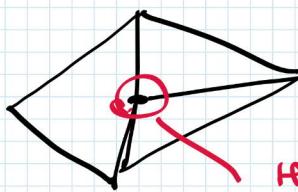
Theorem If $f \in C^0(\bar{\Sigma})$ and it's pieamgle H' , then $f \in H'(H)$
So, continuity guarantees conformity ($V_H \subset V'$).

Each geometrical entity must be the same for all the elements sharing it:

A VERTEX is a vertex for all the elements sharing it.



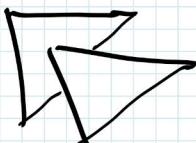
OK



NO

HANGING NODE
This point is a vertex only for some elements

Triangles should not self-intersect or overlap

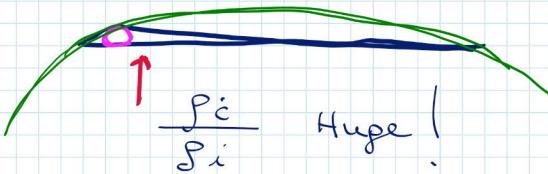
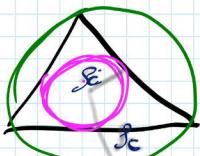


is not acceptable.

REMARK

In some cases we can live with HANGING NODES, but they need a particular care. We do not cover this here.

(4) A good mesh generally prefers "regular" meshes, i.e. triangles with a good "ratio" or quotient between the inscribed and the circumscribed circles.

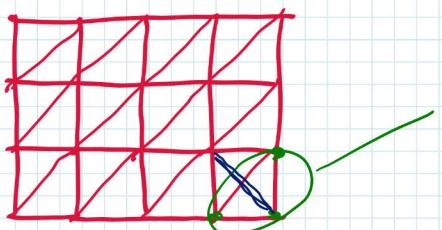


We consider here meshes where

$$\frac{S_c}{S_i} \text{ is bounded}$$

Nevertheless, meshes with extremely skewed elements are sometimes useful (specific theory required).

(5) It's good avoiding all vertices of a triangle on the boundary



having 3 nodes on the boundary, the solution here is completely determined by the B.C., no equations is solved instead \Rightarrow SWAP THE DIAGONAL (Blue)

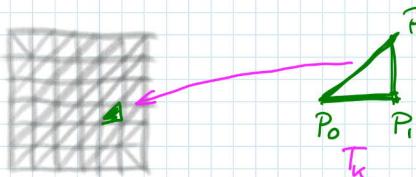
Linear FE in 2D

If we use a triangular mesh with the features mentioned above, we can postulate the construction of a solution in the space:

$$V_h = \{ v_h \in C^0(\bar{\Omega}) : v_h|_{T_k} \in P^r \}$$

where T_k is the generic triangle of the mesh

Here, we are assuming that the mesh is covering exactly the domain ($\Omega_h = \Omega$) like in a square or a rectangle:

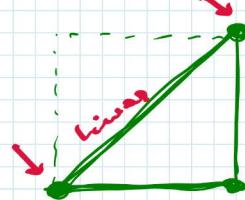


a natural choice that guarantees the continuity of v_h is to use the value in the 3 vertices.

In the generic triangle, we assume the solution to be linear

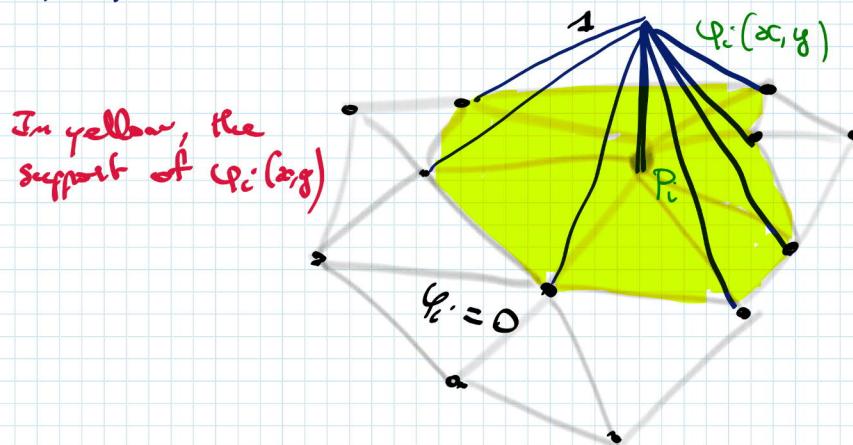
$$u_h(T_k) = ax + by + c$$

We have 3 parameters to identify, so



3 values of u_h in the vertices identify the 3 coefficients a, b, c
A function in the neighbor triangle with the same values at the vertices is automatically continuous.

In a modern perspective, we are constructing piecewise polynomials as illustrated below:

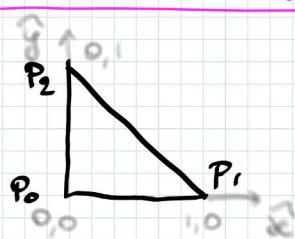


$$u_h = \sum c_j \varphi_j(x, y)$$

with this choice of basis functions: $c_j = u_h(P_j)$

Formally, the solution of the problem is like in the 1D case. However, we have some technical details to cover.

The Reference Finite Element



The factual representation of the basis functions is in two steps

- representation on the reference element
- map current \leftrightarrow reference

Representation on the Reference FE

We have to construct 3 linear functions :

$$\left. \begin{array}{l} \hat{\varphi}_0(P_0) = 1 \\ \hat{\varphi}_0(P_1) = \hat{\varphi}_0(P_2) = 0 \end{array} \right\}$$

Vanishing on the nodes \hat{P}_1 and \hat{P}_2 , and being linear, clearly $\hat{\varphi}_0$ vanishes on the entire segment $\hat{P}_1 \hat{P}_2$.

This line has equation : $\hat{x} + \hat{y} = 1$

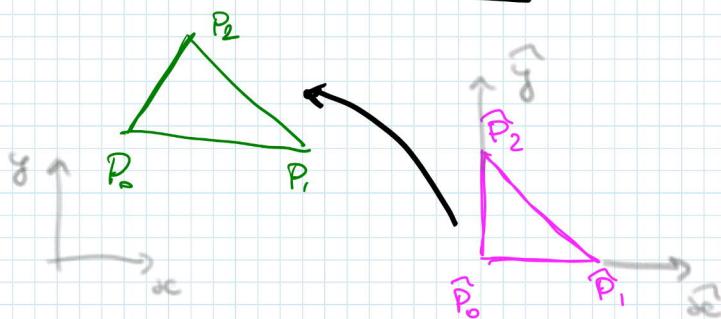
$$\Rightarrow \hat{\varphi}_0(\hat{x}, \hat{y}) = 1 - \hat{x} - \hat{y}$$

(vanishes on \hat{P}_1, \hat{P}_2 automatically)

Similarly : $\hat{\varphi}_1(\hat{x}, \hat{y}) = \hat{x}$, $\hat{\varphi}_2(\hat{x}, \hat{y}) = \hat{y}$

(Btw, notice that $\hat{\varphi}_0 + \hat{\varphi}_1 + \hat{\varphi}_2 = 1$. The same holds in 1D and 3D
 \Rightarrow partition of unity property)

Map Current \leftarrow Reference



The map consists of two functions

$$\left. \begin{array}{l} x(\hat{x}, \hat{y}) \\ y(\hat{x}, \hat{y}) \end{array} \right\}$$

s.t. $\left\{ \begin{array}{l} x(0,0) = x_0 \\ y(0,0) = y_0 \\ x(1,0) = x_1 \\ y(1,0) = y_1 \\ x(0,1) = x_2 \\ y(0,1) = y_2 \end{array} \right.$

We have 3×2 conditions, if $x(\hat{x}, \hat{y})$ and $y(\hat{x}, \hat{y})$ are linear functions, we have 3×2 parameters, so we have a perfect match!!!

The explicit solution is easily constructed by using the same $\hat{\varphi}_0, \hat{\varphi}_1, \hat{\varphi}_2$ of the reference finite element:

$$x = x_0 \hat{\varphi}_0(\hat{x}, \hat{y}) + x_1 \hat{\varphi}_1(\hat{x}, \hat{y}) + x_2 \hat{\varphi}_2(\hat{x}, \hat{y})$$

$$y = y_0 \hat{\varphi}_0(\hat{x}, \hat{y}) + y_1 \hat{\varphi}_1(\hat{x}, \hat{y}) + y_2 \hat{\varphi}_2(\hat{x}, \hat{y})$$

This is an AFFINE map, mapping triangles in Triangle. Finite elements using this map are called affine.

Notice that the jacobian matrix of this transform is

$$J = \begin{bmatrix} \frac{\partial x}{\partial \xi} & \frac{\partial y}{\partial \xi} \\ \frac{\partial x}{\partial \eta} & \frac{\partial y}{\partial \eta} \end{bmatrix} \quad |J| = \frac{\partial x}{\partial \xi} \frac{\partial y}{\partial \eta} - \frac{\partial x}{\partial \eta} \frac{\partial y}{\partial \xi}$$

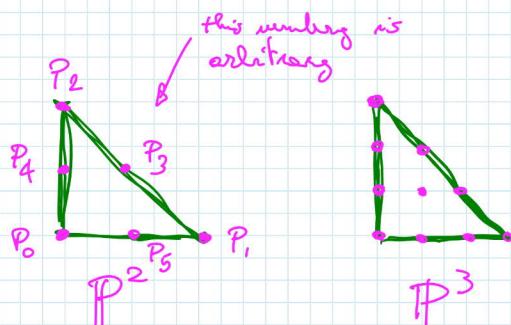
and in general we need to keep in mind that:

$$\nabla u_{\xi\eta} = \begin{bmatrix} \frac{\partial u}{\partial \xi} \\ \frac{\partial u}{\partial \eta} \end{bmatrix} = \begin{bmatrix} \frac{\partial u}{\partial x} \frac{\partial x}{\partial \xi} + \frac{\partial u}{\partial y} \frac{\partial y}{\partial \xi} \\ \frac{\partial u}{\partial x} \frac{\partial x}{\partial \eta} + \frac{\partial u}{\partial y} \frac{\partial y}{\partial \eta} \end{bmatrix} = \begin{bmatrix} \frac{\partial u}{\partial x} & \frac{\partial u}{\partial y} \\ \frac{\partial x}{\partial \xi} & \frac{\partial y}{\partial \xi} \\ \frac{\partial x}{\partial \eta} & \frac{\partial y}{\partial \eta} \end{bmatrix} \nabla_{x,y} u$$

Affine Finite Elements of Order > 1 (2D)

We can obviously introduce finite elements of degree 2 and 3 on elements obtained by an affine map.

On the reference element



\bullet = Lagrangian degrees of freedom

$$u = a x^3 + b x^2 y + c x y^2 + d y^3 + e x^2 + f y^2 + g x y + h x + i y + j$$

6 parameters

$$\hat{\varphi}_0 = 2(-\hat{x}-\hat{y})\left(\frac{1}{2}-\hat{x}-\hat{y}\right)$$

$$\hat{\varphi}_1 = \hat{x}\left(\hat{x}-\frac{1}{2}\right)^2$$

$$\hat{\varphi}_2 = \hat{y}\left(\hat{y}-\frac{1}{2}\right)^2$$

$$\hat{\varphi}_3 = 4\hat{x}\hat{y}$$

$$\hat{\varphi}_4 = 4\hat{y}(1-\hat{x}-\hat{y})$$

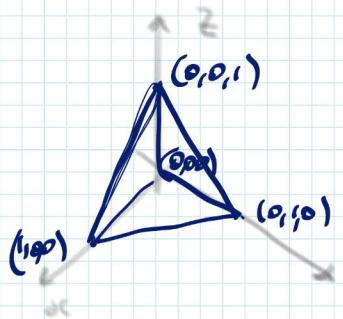
$$\hat{\varphi}_5 = 4\hat{x}(1-\hat{x}-\hat{y})$$

10 parameters

4 dof on the boundary + 1 bubble.

Affine Finite Elements in 3D

In 3D we work directly on the reference element, the unit tetrahedron

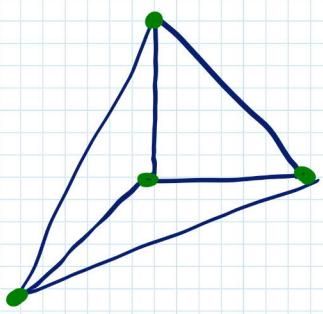


4 vertices

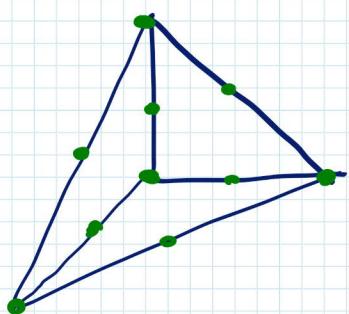
linear functions are in the form :

$$ex+fy+cz+d \Rightarrow 4 \text{ parameters}$$

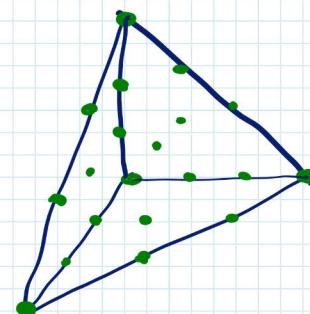
The 4 vertices are used for d.o.f. in the linear case



3 vertices

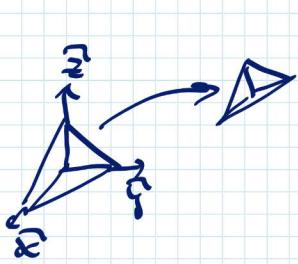


4 vertices +
6 edges = 18 dof



4 vertices + 12 edges +
4 faces = 24 dof

Also in this case the mesh current \leftarrow reference is done by using an affine map with the P^1 basis functions:



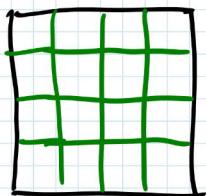
$$\left\{ \begin{array}{l} x = x_0 \hat{\varphi}_0(\bar{x}, \bar{y}, \bar{z}) + x_1 \hat{\varphi}_1(\bar{x}, \bar{y}, \bar{z}) + x_2 \hat{\varphi}_2(\bar{x}, \bar{y}, \bar{z}) \\ \quad + x_3 \hat{\varphi}_3(\bar{x}, \bar{y}, \bar{z}) \\ y = y_0 \hat{\varphi}_0 + y_1 \hat{\varphi}_1 + y_2 \hat{\varphi}_2 + y_3 \hat{\varphi}_3 \\ z = z_0 \hat{\varphi}_0 + z_1 \hat{\varphi}_1 + z_2 \hat{\varphi}_2 + z_3 \hat{\varphi}_3 \end{array} \right.$$

Remarks

(1) Quadrilateral finite elements (and more)

We can consider also different basic geometries, like QUADRILATERALS.

For instance, if we cover a domain with squares:



we need to consider different functions:

(1) Bilinear functions: $axy + bxc + cya + d$ \Rightarrow (linear in each variable)

4 d.o.f.
if
4 equations = 4 vertices

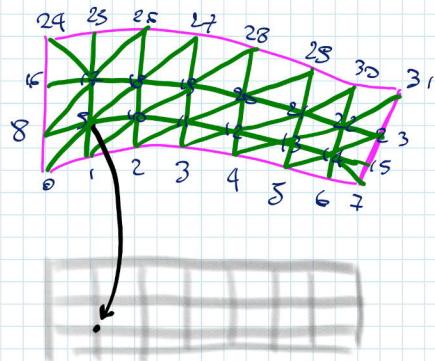
(2) Biquadratic or more functions.

We do not go in details here: bilinear f.e. may have good properties, but the meshing is less versatile.

Recently, Brezzi et al. introduced the so-called **METRIC FINITE DIFFERENCES** that use smart finite elements working on polygons (or, in general, polytopes).

(2) Structured and Unstructured Meshes

When a mesh has a **STRONG CARTESIAN FOOTPRINT** (like in Finite Differences) we can easily infer a relation between the position and the numbering of the nodes/elements.



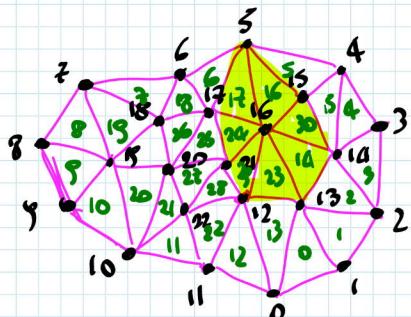
In this case (see the figure) we say that the mesh is **STRUCTURED**.

Having a rule between position and numbering can make the coding more efficient.

However, in most of the cases this is not possible.

In general, the mesh is then called **UNSTRUCTURED**.

In this case, the connection between **NUMBERING & COORDINATES** is all based on the tables in the mesh file.



IMPORTANT

A typical feature of FEM is that the stiffness matrix is **SPARSE**.

With unstructured meshes the sparsity pattern is however disordered.

Nevertheless, it is completely determined by the mesh numbering. This means that we know the pattern of the matrix directly from the mesh.

For instance, in the mesh above we have 22 nodes, 12 of them on the boundary.

For instance, node 16 "sees" the nodes:

16

15 - 5 - 17 - 21 - 12 - 13 - 14 (row 16 has nonzero entries in these columns)

15

18 - 7 - 8 - 9 - 10 - 20

As you can see, the pattern (command `SPY`) will be with no structure.

Isoparametric Finite Elements

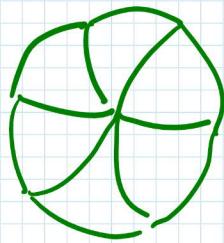
So far, we assumed that $\omega_e = \omega$. However, in general this is not true:
 a circle, for instance, cannot be covered exactly by linear triangles



This introduces a geometrical error that affects the actual convergence rates

\Rightarrow if you are interested,
 ERN, GVERMOND, Theory & Practice of FE, SPRINGER.

This can be avoided using "geometric finite elements":



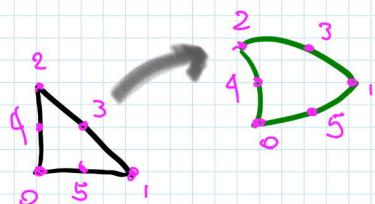
How can we generate those elements?

Well, we can still start from the usual Reference Element.

However, we postulate an either geometrical cubic map that can be described by the same basis functions at the finite elements.

EXAMPLE

Quadratic Elements:



$$\begin{cases} x = \alpha_0 \hat{\varphi}_0 + \alpha_1 \hat{\varphi}_1 + \alpha_2 \hat{\varphi}_2 + \alpha_3 \hat{\varphi}_3 + \alpha_4 \hat{\varphi}_4 + \alpha_5 \hat{\varphi}_5 \\ y = \beta_0 \hat{\varphi}_0 + \dots + \beta_5 \hat{\varphi}_5 \end{cases}$$

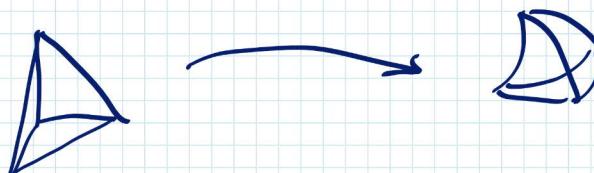
If we use \mathbb{P}^2 with quadratic elements, we have ISOPARAMETRIC ELEMENTS (iso = equal, same; parametric = # of parameters)

\Rightarrow we use the same basis functions for the map and the solution u_h .

(Linear Affine Finite Elements are isoparametric too).

Using these elements, the geometrical error for ω_e being $\neq \omega$ is significantly reduced.

We can obviously consider isoparametric finite elements in 3D too



Convergence Rate and Other Theoretical Aspects in 2+D.

In spite of an increased complexity, many theoretical results can be used in 1D hold in multiple dimensions.

CONVERGENCE RATE :

$$\|u - u_h\|_{H^r} \leq Ch^r \|u\|_{H^{s+1}} \quad r = \min(s, p)$$

for $u \in H^{s+1}(\Omega)$ and with FE at order p .

This can be proved by combining two results :

Daug-Hans Lemma] with regular mesh

Bramble-Hilbert Theorem] $\left(\frac{S_c}{S_i} \leq T < \infty \right)$

and $S_c = S_i$.

CONDITION NUMBER

$$\text{cond}(A) \propto h^{-2} \quad \text{also in 2+D.}$$

Remark : ISOGEOOMETRIC ANALYSIS (IGA)

Recently (≈ 2005) T. Hughes introduced the so-called ISOGEOOMETRIC ANALYSIS

The basic idea is to use NURBS instead of regular polynomials in the elements. NURBS stands for Non-Uniform Rational B-Splines.

These are (piecewise) regular functions much more regular and completed that can attain high convergence rates.

The big advantage of IGA is that NURBS are the basis of CAD software. With IGA the combination of CAD (Computer Aided Design) and CAE (Computer Aided Engineering) is very immediate and having the same approach of isoparametric finite elements, the geometrical accuracy of the domain Ω is very high.

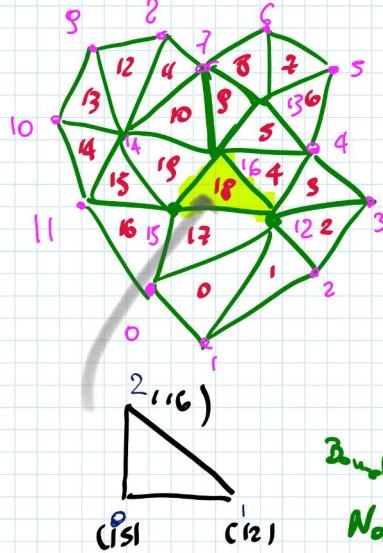
Some Details on the Implementation

As we noticed, the assembly of the stiffness matrix is done elementwise, with a "think globally, act locally" approach.

We first assemble the local matrix on the single element: the local matrix is small and full.

Then, we map the contribution of the local matrix to the global one (exploiting the additivity of the integral operator), such that for instance

$$\int_{\Omega} u_h v_h = \sum_k \int_{T_h} u_h v_h$$

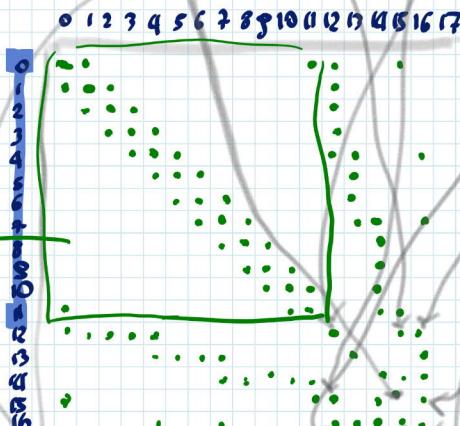


For instance, for linear FEM in 2D, the local matrix is 3×3 :

$$A_{loc} = \begin{bmatrix} a(\hat{\varphi}_0, \hat{\varphi}_0) & a(\hat{\varphi}_1, \hat{\varphi}_0) & a(\hat{\varphi}_2, \hat{\varphi}_0) \\ a(\hat{\varphi}_0, \hat{\varphi}_1) & a(\hat{\varphi}_1, \hat{\varphi}_1) & a(\hat{\varphi}_2, \hat{\varphi}_1) \\ a(\hat{\varphi}_0, \hat{\varphi}_2) & a(\hat{\varphi}_1, \hat{\varphi}_2) & a(\hat{\varphi}_2, \hat{\varphi}_2) \end{bmatrix}$$

Element 18 (yellow)
15
12
16

If all the nodes are Dirichlet,
the matrix is replaced
by the identity matrix



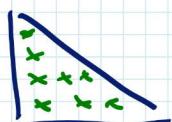
We address the entries by using the mesh file as described earlier.

Element 18 : 15 12 16
↑ ↑ ↗

rows and columns
touched in Δ when processing
element 18.

Computing the Integrals

We used Gaussian or Gauss-Lobatto formulas



$$\int_T f \approx \sum w_i f(x_i, y_i)$$

weights, nodes
tabulated -

Since we do not want errors in the integrals, we try to use a formula with many nodes (= high degree of exactness).

However, be careful! The computational costs rapidly increases.

Boundary Conditions.

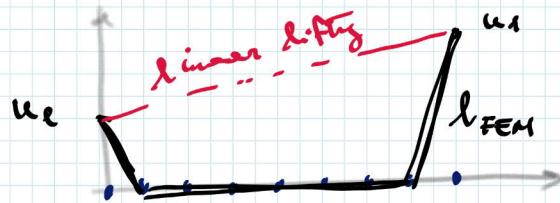
Let's visualize this in 1D, just for simplicity.

When we identify the ^{Doubtless} boundary nodes and set:

$$A(\text{label-D}, \text{label-D}) = 1$$

$$\Delta(\text{label-D}, \text{off-diagonal}) = 0$$

is equivalent to introducing the following lifting of the state:



you can test yourself
the code with an
exercise .

Homework: install Fenics

Week 9

Quateroni
Chap. 12

Advection-Diffusion-Reaction with FEM

We have noticed with FD that sometimes the real numbers require restrictions to the numerical discretization.

In particular, with the problem:

$$-\mu u'' + \beta u' = 0 \quad x \in (0,1) \quad (*)$$

$$u(0) = 0 \quad u(1) = 1$$

we found that the classical centered 2nd order difference scheme:

$$-\mu \frac{u_{i+1} - 2u_i + u_{i-1}}{h^2} + \beta \frac{u_{i+1} - u_{i-1}}{2h} = 0$$

$$u_0 = 0 \quad u_{N_h} = 1 \quad (Nh=1)$$

is oscillating for $\frac{|\beta| h}{2\mu} > 1$

while the solution is monotone.

A natural question is: what happens with FEM?

Do we have the same problem?

What are the solutions?

What to do in 2D?

In this week, we will answer to these questions.

Advection Dominated Problems in 1D

To start with, we can check what we obtain if we discretize the same problem (*) with linear finite elements.

The variational form reads:

$$\mu \int_0^1 u'_h \varphi'_i + \beta \int_0^1 u'_h \varphi_i = 0 \quad \varphi_i(0) = \varphi_i(1) = 0$$

The entries for the first integral are:

$$\mu \int_0^1 \varphi'_j \cdot \varphi'_i = \begin{cases} \frac{2}{h} \mu & i=j \\ -\frac{1}{h} \mu & j=i \pm 1 \\ 0 & \text{if } i,j \neq 1 \end{cases}$$

For the second one:

$$\beta \int_0^1 \varphi_i \varphi_j' = \begin{cases} 0 & j=i \text{ and } \text{for } j \neq i \pm 1 \\ \beta & j=i+1 \\ -\beta & j=i-1 \end{cases}$$

So, at the end we get the scheme:

$$-\mu \frac{u_{i+1} - 2u_i + u_{i-1}}{h} + \beta(u_{i+1} - u_{i-1}) = 0$$

$$u_0 = 0 \quad u_{N+1} = 1$$

But THIS IS EXACTLY THE FD SCHEME, just multiplied by μ .

\Rightarrow Also with FEM we should expect oscillations for $\frac{\beta h}{2\mu} > 1$.

We already knew that upwind is a remedy, but the original physics (\approx go toward the upwind direction) doesn't apply to the variational formulation. In fact, that interpretation works only for FD.

To apply the concept to FEM, we resort to the ARTIFICIAL VISCOSITY.

We found that Upwind in 1D is equivalent to an increment of the coefficient μ , so that:

$$\mu^* = \mu(1+\alpha)$$

With FE, we can therefore suppress the oscillation by solving the following problem:

$$\mu^* \int_0^1 u_h v_h + \frac{1+\alpha}{2} \int_0^1 u_h v_h' + \beta \int_0^1 u_h' v_h = 0 \quad u_h(0) = 0, \quad u_h(1) = 1, \quad \forall v_h \in H_0^1(0,1).$$

For the FD problem, we know that this increment of the viscosity was obtained by reducing the accuracy and that we need a second-order upwind scheme to preserve the overall 2nd-order of the FD scheme.

What can we say for the FE perspective (i.e., error analysis)?

A first remark: the real term is a "perturbation" of the real problem we want to solve. However, it's a legal perturbation as long as it vanishes with h .

Let's try to formalize the concept.

The original problem reads:

$$\text{find } u \in V \text{ s.t. } \alpha(u, v) = f(v) \quad \forall v \in V$$

$$\text{where } \alpha(u, v) = \mu \int_0^1 u v' + \beta \int_0^1 u' v \quad \text{and } f(v) = 0.$$

The numerical problem with upwind reads:

$$\text{Find } u_h \in V_h \text{ s.t. } \alpha_h(u_h, v_h) = f(v_h) \quad \forall v_h \in V_h$$

$$\text{where } \alpha_h(u_h, v_h) = \alpha(u_h, v_h) + \frac{\beta h}{2} \int_0^1 u'_h v'_h dh$$

$$(\text{so that } \alpha_h(u_h, v_h) \rightarrow \alpha(u_h, v_h) \text{ for } h \rightarrow 0).$$

A first important consequence is that THE CÉA LEMMA DOES NOT HOLD ANY LONGER.

$$\text{While } \alpha(u, v_h) = \alpha(u_h, v_h), \text{ we have:}$$

$$\alpha_h(u, v_h) \neq f(v_h), \Rightarrow \alpha_h(u - u_h, v_h) \neq 0.$$

\Rightarrow THE UPWIND FE METHOD IS NOT STRANGELY CONSISTENT.

$$\text{However } |\alpha_h(u_h, v_h) - \alpha(u_h, v_h)| = \left| \frac{\beta h}{2} \int_0^1 u'_h v'_h dh \right| \xrightarrow{h \rightarrow 0} 0.$$

$$\text{We conclude that } \alpha_h(u - u_h, v_h) = f(v_h) - f(v_h) - \frac{\beta h}{2} \int_0^1 u'_h v'_h dh = \gamma(h) \xrightarrow{h \rightarrow 0} 0$$

so Upwind is CONSISTENT (not only consistent).

From the Lax-Richtmyer Theorem (convergence = stability + consistency) we know that this is enough for the convergence.

As a matter of fact, it is easy to prove that our upwind scheme is stable:

$$\alpha_h(u_h, v_h) = \alpha(u_h, v_h) + \frac{\beta h}{2} \int_0^1 u'_h v'_h dh \quad \text{is}$$

BILINEAR
CONTINUOUS
COERCIVE

$$\alpha_h(u_h, v_h) = \underbrace{\alpha(u_h, v_h)}_{\text{coercive}} + \underbrace{\frac{\beta h}{2} \int_0^1 u'_h v'_h dh}_{>0} \geq \alpha \|u_h\|_H^2$$

$$\Rightarrow \alpha \|u_h\|_H^2 \leq \alpha_h(u_h, u_h) \leq \|f(u_h)\| \leq M \|u_h\|_H \Rightarrow \|u_h\|_H \leq M/\alpha.$$

OK, we have the convergence, but what about the convergence rate?
We don't have the Céa Lemma.

We need a new theoretical tool.

This is given by the so-called STRANG-LEMMA.

Generalized Galerkin Schemes (GGS)

Let's consider, in general, the following class of schemes for our original problem, called Generalized Galerkin:

$$\text{Find } u_h \in V_h, \text{ s.t. } \alpha_h(u_h, v_h) = f_h(v_h) \quad \forall v_h \in V_h \quad (\star\star)$$

where

$$\alpha_h(u_h, v_h) - \alpha(u_h, v_h) = b_h(u_h, v_h)$$

$$f_h(v_h) - f(v_h) = g_h(u_h, v_h)$$

(for usual this term is 0)

STRANG LEMMA

Let's assume that $a(\cdot, \cdot)$ is bilinear, continuous and coercive

$\mathcal{G}_h(\cdot, \cdot)$ is linear, continuous

(with continuity and coercivity constants independent of h)

Then, the Generalized Galerkin problem is well-posed and \exists Constants C_1, C_2, C_3 s.t.

$$\|u - u_h\|_V \leq \inf_{w_h \in V_h} \left(C_1 \|u - w_h\|_V + C_2 \sup_{v_h \in V_h} \frac{|\mathcal{L}_h(w_h, v_h)|}{\|v_h\|_V} \right) + C_3 \sup_{v_h \in V_h} \frac{|\mathcal{G}_h(v_h)|}{\|v_h\|_V}$$

↑
discretization
error
(as with the
Céa Lemma)

error for the
changes in a_h and \mathcal{G}_h

This Lemma generalizes to CG the Céa Lemma, showing that the total error is the combination of the discretization and the modifications of $a(\cdot, \cdot)$ and $\mathcal{G}(\cdot)$.

In the case of UPWIND, we have

$$\mathcal{L}_h(w_h, v_h) = \frac{\beta h}{2} \int w_h' v_h = \frac{|\mathcal{L}(w_h, v_h)|}{\|v_h\|_V} \leq \alpha \frac{|\beta|}{2} \|w_h\|_V$$

$$|\mathcal{G}_h(v_h)| = 0$$

so we conclude that $\|u - u_h^{\text{UPW}}\|_V \leq Ch$ regardless of the FE degree.

REMARK

The Strang Lemma is useful also for estimating the impact of numerical quadrature on the performance of the method.

In fact, we can write:

$$a(u, v) = \int_0^1 \text{"some terms"}$$

$$a_h(u, v) = \text{Numerical Quadrature ("some terms")}$$

In this respect, $\mathcal{L}_h(\cdot, \cdot)$ and $\mathcal{G}_h(\cdot)$ are the quadrature errors.

If we use Gaussian formulas with an adequate number of nodes, the convergence rate of these terms will be always $> r$ ($= \min(p, \delta)$) of the discretization term. This requires an adequate number of nodes to use, which can be computationally expensive.

However, there are cases where high-order quadrature formulas are not necessary. Or not even convenient.

For instance, for constant coefficients problems with FEM of order q

$$a(\varphi_j, \varphi_i) = \mu \int_0^1 \varphi_j' \varphi_i' + \beta \int_0^1 \varphi_j' \varphi_i + \sigma \int_0^1 \varphi_j \varphi_i$$

↑
 deg. $q-1 \times q-1$
 2q-2

9-1 x 9
 2q-1
 q x q
 2q

With a quadrature with DEGREE OF EXACTNESS $2q$ or more we do not have any error for the quadrature.

This shows how important is to know the concept of Degree of Exactness (the highest polynomial degree integrated exactly by a quadrature formula).

As we said, in general, we can write:

$$\begin{aligned} b_h(u_h, v_h) &= a(u_h, v_h) - \text{Numerical Integration of } a(v_h, v_h) \\ g_h(v_h) &= f(v_h) - " " " f(v_h, v_h) \end{aligned}$$

If the quadrature formula is accurate enough, these terms feature a dependence on h^{q+1} with $q+1 >$ discretization order, so the quadrature error is present but it doesn't affect the convergence order due to the discretization.

This obviously requires a certain number of Gaussian quadrature nodes, with an impact on the computational costs.

High Order Upwind (in 1D) with FE

We noticed that with FD we can get higher order stable schemes by using non-central high order formulas. What can we do with FE, where we do not have a "centred" or "non-centred" choice?

In 1D we can stick to the concept of artificial viscosity: we found that the basic upwind is an artificial viscosity s.t.

$$\mu^* = \mu(1 + \bar{\Phi})$$

The idea is to consider an artificial viscosity:

$$\mu^{**} = \mu(1 + \bar{\Phi}(\bar{\Phi}))$$

where $\bar{\Phi}(\bar{\Phi})$ is a function with the following features:

- ✓ $\tilde{\Phi}(\text{Pe})$ "large enough" (stabilization): $\frac{|\beta| \epsilon}{2\mu \text{Pe}} < 1 \Rightarrow \frac{\text{Pe}}{1 + \tilde{\Phi}(\text{Pe})} < 1 \Rightarrow \tilde{\Phi}(\text{Pe}) - \text{Pe} + 1 > 0$
- ✓ $\lim_{\text{Pe} \rightarrow 0} \tilde{\Phi}(\text{Pe}) = 0$ and possibly $= O(h^p)$ $p > 1$
(accuracy: $(\log(\cdot) \sim O(h^p))$)

A possible function with these features is the so-called SCHAFTER-GUILLÉ function with some features is the so-called SCHAFTER-GUILLÉ stabilization:

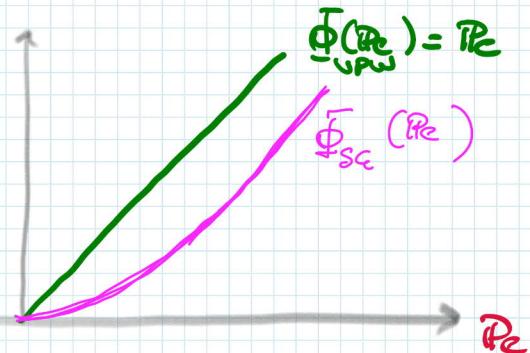
$$\tilde{\Phi}_{SG}(\text{Pe}) = \text{Pe} - 1 + \frac{\text{Pe}}{e^{\text{Pe}-1}} \quad (\text{with } \tilde{\Phi}_{SG}(0) = 0 \text{ to manage the singularity})$$

This function is "smart": clearly

$$\tilde{\Phi}_{SG}(\text{Pe}) - \text{Pe} + 1 = \frac{\text{Pe}}{e^{\text{Pe}-1}} > 0$$

and

$$\lim_{h \rightarrow 0} \tilde{\Phi}_{SG}(\text{Pe}) \sim O(h^2)$$



For high Pe , $\tilde{\Phi}_{SG}(\text{Pe}) \approx \tilde{\Phi}_{UPW}(\text{Pe})$, for h small, it converges to 0 faster than UPW.

SG is so good that it is possible to move kets in the problem:

$$-\mu u' + \beta u' = f$$

for a piecewise constant f , it gives a "modally exact" solution.

Upwind and SG as Petrov-Galerkin methods.

The basic Galerkin method reads:

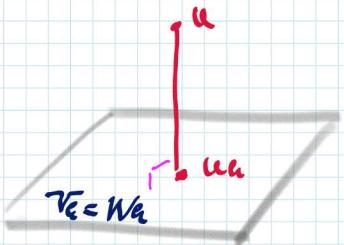
$$\varphi(u_h, v_h) = \mathcal{F}(v_h)$$

with u_h, v_h belonging to the same space.

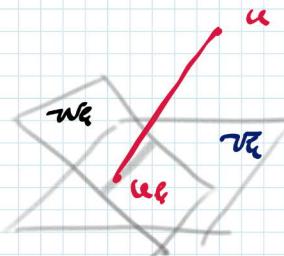
This leads to $\varphi(u - u_h, v_h) = 0$ that is an orthogonal projection.

We can generalize to an oblique projection with:

$$\varphi(u_h, v_h) = \mathcal{F}(v_h) \quad \begin{array}{l} \text{where } u_h \text{ has the same dimension} \\ \text{(to have a square problem) but} \\ \text{it is not } = v_h. \end{array}$$



Orthogonal Projection
(Galerkin)



Oblique Projection
(Petrov-Galerkin)

It is possible to see for LINEAR FINITE ELEMENTS, that u_h and u can be represented as Petrov-Galerkin methods, where

$$V_h = V_h + \{B^\alpha\}$$

where B are "bubble" functions and $\alpha = 1$ for UPW and another "strange" value for SG.

See Quarteroni 12.8.2

Reaction Dominated Problems

Let's consider this problem:

$$\begin{aligned} -\mu u'' + \sigma u &= 0 \\ u(0) &= 0 \quad u(1) = 1 \end{aligned}$$

$$\begin{aligned} \alpha &\in (0, 1) \\ (\mu > 0, \sigma \geq 0) \end{aligned}$$

Should we expect numerical instabilities when $\sigma \gg \mu$ as it happens for advection dominated problems?

The problem is certainly well posed:

$$a(u, v) = \int_0^1 u' v' + \sigma \int_0^1 u v$$

$$\left\{ \begin{array}{ll} \text{Continuous with continuity constant: } M = \max(\sigma, \mu) \\ \text{Coercive with coercivity constant: } \alpha = \min(\sigma, \mu) \end{array} \right.$$

so the numerical approximation is such that

$$\|u - u_h\|_0 \leq \frac{M}{\alpha} \inf_{w \in V_h} \|u - w\|_0$$

If $\sigma \gg \mu$ we have $\frac{M}{\alpha} = \frac{\sigma}{\mu} \gg 1$ and we may expect some instabilities.

Now, let's consider in detail the discretization of the problem.

(1) Exact solution:

$$u'' = \frac{\sigma}{\mu} u$$

$$u = C_1 e^{\sqrt{\frac{\sigma}{\mu}}x} + C_2 e^{-\sqrt{\frac{\sigma}{\mu}}x}$$

$$u(0) = 0 : C_1 = -C_2$$

$$u(1) = 1 : C_1 e^{\sqrt{\frac{\sigma}{\mu}}} + C_2 e^{-\sqrt{\frac{\sigma}{\mu}}} = 1 \quad \left. \begin{array}{l} C_1 = \frac{1}{e^{\sqrt{\frac{\sigma}{\mu}}} - e^{-\sqrt{\frac{\sigma}{\mu}}}} \end{array} \right\}$$

$$u_{ex} = \frac{e^{\sqrt{\frac{\sigma}{\mu}}x} - e^{-\sqrt{\frac{\sigma}{\mu}}x}}{e^{\sqrt{\frac{\sigma}{\mu}}} - e^{-\sqrt{\frac{\sigma}{\mu}}}}$$

2) Finite Difference solution:

$$\left. \begin{array}{l} -\mu \frac{(u_{i+1} - 2u_i + u_{i-1})}{h^2} + \sigma u_i = 0 \\ u_0 = 0 \quad u_{N+1} = 1 \end{array} \right\} \begin{array}{l} u_{i+1} - 2\left(1 + \frac{\sigma h^2}{2\mu}\right) u_i + u_{i-1} = 0 \\ g_{1,2} = \left(1 + \frac{\sigma h^2}{2\mu}\right) \pm \sqrt{\left(1 + \frac{\sigma h^2}{2\mu}\right)^2 - 1} = \\ = 1 + \frac{\sigma h^2}{2\mu} \pm \sqrt{\frac{\sigma h^2}{\mu} + \frac{\sigma^2 h^4}{4\mu^2}} > 0 \quad \forall h \end{array}$$

Since g_1, g_2 are always > 0 , no oscillations!

(3) Linear Finite Elements discretization

We already know that the diffusive term leads to:

$$-\mu u'' \Rightarrow \int_0^1 u'' dx \Rightarrow -\mu \frac{u_{i+1} - 2u_i + u_{i-1}}{h}$$

For the reactive term:

$$\sigma u \Rightarrow \sigma \int_0^1 u v \Rightarrow ?$$

$$\sigma \int_0^1 u^2 = \sigma \int_{x_{i-1}}^{x_i} \frac{(x - x_{i-1})^2}{h^2} + \sigma \int_{x_i}^{x_{i+1}} \frac{(x_{i+1} - x)^2}{h^2} = \\ = \frac{\sigma}{h^2} \frac{1}{3} h^3 + \frac{\sigma}{h^2} \frac{1}{3} h^3 = \frac{2}{3} \sigma h$$

$$\sigma \int_0^1 \varphi_i \varphi_{i-1} = \sigma \int_{x_{i-1}}^{x_i} \frac{1}{h^2} (x_{i-1} - x)(x - x_i) = \dots = \sigma h \frac{1}{6}$$

$$\sigma \int_0^1 \varphi_i \varphi_{i+1} = \dots = \sigma h \frac{1}{6}$$

So, at the end we have:

$$-\mu \frac{u_{i+1} - 2u_i + u_{i-1}}{h} + \frac{\sigma h}{6} (u_{i+1} + 4u_i + u_{i-1}) = 0$$

$$\frac{\mu}{h} \left(\frac{\sigma h^2}{6\mu} - 1 \right) u_{i+1} + \frac{\sigma h}{h} \left(1 + \frac{\sigma h^2}{3\mu} \right) (4 + \frac{\mu}{h} \left(\frac{\sigma h^2}{6\mu} - 1 \right)) u_{i-1} = 0$$

$$u_{i+1} - 2 \begin{bmatrix} 1 + \frac{\sigma h^2}{3\mu} \\ 1 - \frac{\sigma h^2}{6\mu} \end{bmatrix} u_i + u_{i-1} = 0$$

$$\gamma^2 - 2\gamma s + 1 = 0 \quad s_{1,2} = \gamma \pm \sqrt{\gamma^2 - 1}$$

Notice that, if $1 - \frac{\sigma h^2}{6\mu} > 0$, then $\gamma > 0$ and precisely:

$$\gamma = \frac{1 + \frac{\sigma h^2}{3\mu}}{1 - \frac{\sigma h^2}{6\mu}} = \frac{1 - \frac{\sigma h^2}{3\mu} + \frac{\sigma h^2}{3\mu} + \frac{\sigma h^2}{6\mu}}{1 - \frac{\sigma h^2}{6\mu}} > 1 \quad \forall h$$

Then :

$$s_{1,2} = \gamma \pm \sqrt{\gamma^2 - 1} > 0 \quad \forall h$$

In the case $1 - \frac{\sigma h^2}{6\mu} < 0$, then:

$$|\gamma| = \frac{1 + \frac{\sigma h^2}{3\mu}}{\frac{\sigma h^2}{6\mu} - 1} = \frac{\frac{\sigma h^2}{6\mu} - 1 + \frac{\sigma h^2}{6\mu} + 2}{\frac{\sigma h^2}{6\mu} - 1} = 1 + \frac{\frac{\sigma h^2}{6\mu} + 2}{\frac{\sigma h^2}{6\mu} - 1} > 1$$

Then

$$s_{1,2} = -|\gamma| \pm \sqrt{\gamma^2 - 1} \quad \text{are real but one is negative}$$

\Rightarrow We have oscillations!

Error Finite Elements present oscillations if $\frac{\sigma h^2}{6\mu} > 1$

We can introduce the Reactive Peclét number and state that to avoid oscillations we need:

$$Pe_R = \frac{\sigma h^2}{6\mu} < 1 \Rightarrow R < \sqrt{\frac{6\mu}{\sigma}}$$

Differently from FD, FE do suffer in the case $\frac{\mu}{G} \ll 1$ (as expected by our analysis).

Questions: do we have a workaround?
What the FD can teach us?

The difference with the FD method is in the term once that

with FD leads to σu_i , with FE to $\frac{\sigma h}{6} (u_{i+1} + 4u_i + u_{i-1})$.

Notice that the reactive term is given by the matrix

$$M_{ij} = \int_0^1 \varphi_i \varphi_j \quad \text{that is called MASS MATRIX}$$

$M_{ij} = \begin{cases} \frac{\sigma h}{6} & j=i \\ \frac{2\sigma h}{3} & j=i+1 \\ 0 & \text{otherwise} \end{cases}$

(we will see why)

Also, we can notice that the integral $f = \varphi_i \varphi_j$ is a quadratic term so a quadrature formula with Degree of Exactness ≥ 2 is enough to compute it exactly. Simpson for instance:

$$(\text{Simpson}) \int_a^b f = \frac{\sigma h}{6} \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right)$$

In fact:

$$\int_{x_i}^{x_{i+1}} \varphi_i \varphi_{i+1} \underset{\text{Simpson}}{\approx} = \frac{\sigma h}{6} \left(f(x_i) + 4f\left(\frac{x_i+x_{i+1}}{2}\right) + f(x_{i+1}) \right) =$$

$$= \frac{\sigma h}{6} \frac{1}{h^2} \frac{h^2}{4} = \frac{\sigma h}{6} \quad \text{EXACT},$$

$$\int_{x_{i-1}}^{x_i} \varphi_i^2 = \int_{x_{i-1}}^{x_i} \varphi_i^2 \text{ to } \int_{x_i}^{x_{i+1}} \varphi_i^2 \underset{\text{(Simpson)}}{\approx} = \frac{\sigma h}{6} \left(\cancel{\varphi_i^2(x_{i-1})} + 4\varphi_i^2\left(\frac{x_i+x_{i+1}}{2}\right) + \cancel{2\varphi_i^2(x_{i+1})} \right. \\ \left. + 4\varphi_i^2\left(\frac{x_i+x_{i+1}}{2}\right) + \cancel{\varphi_i^2(x_{i+1})} \right) = \\ = \frac{\sigma h}{6} \left(4 \frac{1}{4} + 2 + 4 \frac{1}{4} \right) = \frac{4\sigma h}{6} \quad \text{EXACT}$$

On the contrary, if we use a formula like the Trapezoidal one, we get a quadrature error:

$$(\text{T}) \int_a^b f = \frac{b-a}{2} (f(a) + f(b)) \quad (\Delta x = 1)$$

$$\text{T} \int_{x_{i-1}}^{x_i} \varphi_i \varphi_{i+1} \underset{\text{(T)}}{\approx} 0 \quad (\varphi_i(x_{i+1}) = 0 \quad \text{and} \quad \varphi_{i+1}(x_i) = 0)$$

$$\sigma \int_{x_{i-1}}^{x_i} \varphi_i^2 = \sigma \int_{x_{i-1}}^{x_i} \varphi_i^2 + \sigma \int_{x_i}^{x_{i+1}} \varphi_i^2 \underset{\text{(T)}}{\approx} \frac{\sigma h}{2} \left(\cancel{\varphi_i^2(x_{i-1})} + 2\boxed{\varphi_i^2(x_i) + \cancel{\varphi_i^2(x_{i+1})}} \right) = \\ = \sigma h$$

So, if we use the Trapezoidal Rule for FE (linear) we obtain:

$$-\frac{\mu}{h} (u_{i+1} - 2u_i + u_{i-1}) + \sigma h u_i = 0$$

that is

$$h \underbrace{\left(-\frac{\mu}{h^2} (u_{i+1} - 2u_i + u_{i-1}) + \sigma u_i \right)}_{\text{FD scheme!!!}} = 0$$

For this problem : Trapezoidal + Linear FE $\equiv FD$!!

How can we analyze this case?

We can look at the Trapezoidal integration as a perturbation of the bilinear form $a(\cdot, \cdot)$.

Keep in mind that, in general:

$$\begin{aligned} \text{SIMPLE TRAPEZOIDAL: } & \left| \int_a^b f - \frac{(b-a)}{2} (f(a) + f(b)) \right| = \\ &= \left| \int_a^b f(x) + \int_a^b f'(x)(x-a) + \int_a^b f''(x) \frac{(x-a)^2}{2} - \frac{(b-a)}{2} (f(a) + f(b)) \right| = \\ &= \left| f(a)(b-a) + f'(a) \frac{(b-a)^2}{2} - \frac{(b-a)}{2} (f(a) + f(b) + f'(a)(b-a) + f''(a) \frac{(b-a)^2}{2} + \int_a^b f'''(x) \frac{(x-a)^2}{2}) \right| \\ &\sim O((b-a)^3) \end{aligned}$$

$$\text{COMPOSITE TRAPEZOIDAL: } \left| \int_a^b f - \sum_{i=1}^{N-1} \int_{x_i}^{x_{i+1}} f \right| \sim O(h^2)$$

In this way, we can state:

$$u(u, v) = \underbrace{\text{Trapezoidal}(u(u, v))}_{u_h(u, v)} + O(h^2)$$

When we use the trapezoidal Rule, we fall into the class of Generalized Galerkin method with l_h and g_h scaling with $O(h^2)$.

Then we have from the Strang Lemma:

$$\|u - u_h\|_V \leq \underbrace{C_1 h}_{\text{discretization}} + \underbrace{C_2 h^2 + C_3 h^2}_{\text{quadrature}} \quad (p=1)$$

Our method will still be first order in the L^1 -norm.

The behavior in the L^2 -norm will be tested in our numerical session.

Linear FE + Trapezoidal Rule is called MASS LUMPING :

in fact, it is like summing up the off-diagonal elements to the diagonal

elements of the mass-matrix:

$$M_{FE} = \left[\begin{array}{ccc} \sigma \frac{R}{6} & \xrightarrow{+} & \sigma \frac{4}{6} h \xleftarrow{+} \sigma \frac{h}{6} \\ \textcircled{1} & & \textcircled{2} \end{array} \right] \Rightarrow M_C = \left[\begin{array}{cc} \sigma h & \textcircled{1} \\ \textcircled{2} & \sigma h \end{array} \right]$$

What happens with FE of order 2+?

The mass-lumping doesn't work in the same way: trapezoidal rule would lead to a singular mass matrix.

However, mass-lumping strategies are available \Rightarrow see Question Book.

WEEK 10

Advection-Diffusion-Reaction Problems in 2+ dimensions.

Let's consider this problem:

$$-\mu \Delta u + \beta \cdot \nabla u \quad (\text{con}) = f$$

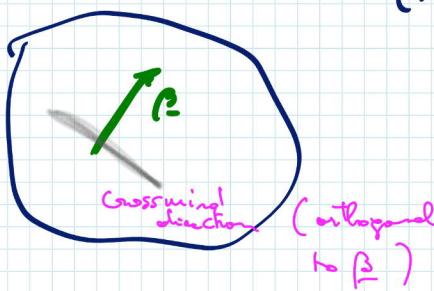
$$u(\partial\Omega) = \text{state}.$$

We know already that:

(1) for $\|\beta\| \gg \mu$ we should expect instabilities

(then the 1D case and the FD experience)

(2) the action of a possible numerical viscosity should not go in all the directions, as only the streamline direction (identified by β) requires special treatments.



Along the crosswind direction, we do not need additional viscosity.

How can we organize the "directional viscosity" in the variational formulation?

(To be honest, I think it's easier with the variational formulation than with the strong one).

Follow: $-\nabla \cdot (\mu \nabla u + \beta u) \quad (\text{con}) = f$

The additional viscosity consists of adding the (isotropic) term:

$$-\nabla \cdot \left(\frac{\|\beta\| h}{2} \nabla u \right)$$

but in fact we need a stabilization only along β , so we should consider to modify the term:

$$-\nabla \cdot \left(\delta h \underbrace{(\beta \cdot \nabla u) \frac{\beta}{\|\beta\|}}_{\text{stabilization along } \beta} \right)$$

$\frac{\beta}{\|\beta\|}$ = unit vector in the direction of β .

δ is an arbitrary parameter to be tuned numerically.

In the variational form, this leads to:

$$\mu \int_{\Omega} \nabla u \cdot \nabla v + \int_{\Omega} (\beta \cdot \nabla u) v + \boxed{\delta h \int_{\Omega} (\beta \cdot \nabla u) \beta \cdot \nabla v}$$

→ This is the term actually stabilized (regardless of the B.C.) when we use $(\delta h \rightarrow)$

Streamline Diffusion (SD) method.

In this way, we resort to a generalized Galerkin method

$$\alpha(u_h, v_h) = \beta(v_h)$$

$$\text{where } \alpha_h(u_h, v_h) = \alpha(u_h, v_h) + \frac{\delta h}{\|\beta\|} \int_{\Omega} (\beta \cdot \nabla u_h)(\beta \cdot \nabla v_h)$$

$$\text{To be explicit: } \frac{\delta h}{\|\beta\|} \int_{\Omega} (\beta \cdot \nabla u_h)(\beta \cdot \nabla v_h) =$$

$$= \frac{\delta h}{\|\beta\|} \int_{\Omega} \left(\beta_1 \frac{\partial u_h}{\partial x_1} + \beta_2 \frac{\partial u_h}{\partial x_2} \right) \left(\beta_1 \frac{\partial v_h}{\partial x_1} + \beta_2 \frac{\partial v_h}{\partial x_2} \right) =$$

$$= \frac{\delta h}{\|\beta\|} \int_{\Omega} \beta_1 \frac{\partial u_h}{\partial x_1} \frac{\partial v_h}{\partial x_1} + \beta_1 \beta_2 \left(\frac{\partial u_h}{\partial x_1} \frac{\partial v_h}{\partial x_2} + \frac{\partial u_h}{\partial x_2} \frac{\partial v_h}{\partial x_1} \right) + \beta_2^2 \frac{\partial u_h}{\partial x_2} \frac{\partial v_h}{\partial x_2}$$

Notice that if $\beta_2 = 0$, we get exactly the stabilization we expect only along x_1 (and the other way around for $\beta_1 = 0$).

The Streamline Diffusion method is very popular, it loses the strong consistency of Galerkin and it can be analyzed by the Strong Frame.

We find straightforwardly that the method is 1st order regardless of the polynomial degree of finite elements.

On the other hand, if $\alpha(u, v)$ is coercive,

$$\alpha_h(u, v) = \alpha(u, v) + \frac{\delta h}{\|\beta\|} \int_{\Omega} (\beta \cdot \nabla u)(\beta \cdot \nabla v)$$

(with $\delta > 0$) is trivially (more) coercive, with a coercivity constant

$$\alpha_{SD} = \alpha + C\delta h$$

so that by an appropriate tuning of δ we can avoid spurious oscillations.

The natural question however is: can we obtain a stable high order (> 1) method?

Instead of tuning a different coefficient δ (as with Schatzfetter-Gummel) we change paradigm:

CAN WE FIND A STRONGLY CONSISTENT STABLE METHOD?

To answer this question, we need some preliminary definitions.

Symmetric and Skewsymmetric parts of an operator.

We already know that if a bilinear form is such that

$$Q(u, v) = Q(v, u) \quad \forall u, v \in V$$

we say that it is Symmetric-

We also say that if a bilinear form is s.s.

$$Q(u, v) = -Q(v, u)$$

we say that it is skew-symmetric

In linear algebra, similarly, we have symmetric and skew-symmetric matrices:

$$A = A^T \quad (\text{symmetric})$$

$$A = -A^T \quad (\text{skew-symmetric})$$

Now, in linear algebra, for a generic matrix A , we have:

$$A = \underbrace{\frac{1}{2}(A + A^T)}_{\text{symmetric}} + \underbrace{\frac{1}{2}(A - A^T)}_{\text{skew-symmetric}}$$

Do we have the same for bilinear forms?

Let's work on the strong formulation:

$$\underbrace{-\mu \Delta u + \beta \cdot \nabla u}_{\int_{\Omega} u \nabla u \cdot \beta \nu} + \underbrace{\alpha u}_{(\text{symmetric})} = f$$

\downarrow
 $\int_{\Omega} u \nabla u \cdot \beta \nu$
 (symmetric)

\downarrow
 $\int_{\Omega} u v$
 (symmetric)

$$\underbrace{\int_{\Omega} \beta \cdot \nabla u v}_{Q_C(u, v)} = \int_{\Omega} \beta \cdot \nabla u v - \int_{\Omega} u \nabla \cdot (\beta v) = \text{Boundary term}$$

$$- \int_{\Omega} (\nabla \cdot \beta) u v - \int_{\Omega} (\beta \cdot \nabla v) u$$

\downarrow
 Symmetric

$$Q_C(v, u)$$

So, from here and dropping the boundary term:

$$\underbrace{\int_{\Omega} (\beta \cdot \nabla u) v + \frac{1}{2} \int_{\Omega} (\nabla \cdot \beta) u v}_{\text{Skew-Symmetric}} = - \int_{\Omega} (\beta \cdot \nabla v) u - \frac{1}{2} \int_{\Omega} (\beta \cdot \beta) u v$$

$$\Rightarrow \int_{\Omega} (\underline{\beta} \cdot \nabla u) v = \underbrace{\int_{\Omega} (\underline{\beta} \cdot \nabla u) v + \frac{1}{2} \int_{\Omega} (\nabla \cdot \underline{\beta}) uv}_{\text{Skewsymmetric}} - \underbrace{\frac{1}{2} \int_{\Omega} (\underline{\beta} \cdot \underline{\beta}) uv}_{\text{Symmetric}}$$

REMARK

In the book of Quarteroni, you find the convective term in the form:

$$\int_{\Omega} \nabla \cdot (\underline{\beta} u) v$$

In this case:

$$\int_{\Omega} \nabla \cdot (\underline{\beta} u) v = \int_{\partial \Omega} \underline{\beta} \cdot \underline{n} uv - \int_{\Omega} (\underline{\beta} \cdot \nabla v) u = \text{Boundary Term} - \int_{\Omega} \nabla \cdot (\underline{\beta} u) v + \int_{\Omega} (\nabla \cdot \underline{\beta}) uv$$

Again, dropping the boundary term:

$$\int_{\Omega} \nabla \cdot (\underline{\beta} u) v - \frac{1}{2} \int_{\Omega} (\nabla \cdot \underline{\beta}) uv = - \int_{\Omega} \nabla \cdot (\underline{\beta} v) u + \frac{1}{2} \int_{\Omega} (\underline{\beta} \cdot \underline{\beta}) uv$$

Skewsymmetric :

$$\int_{\Omega} \nabla \cdot (\underline{\beta} u) v = \underbrace{\int_{\Omega} \nabla \cdot (\underline{\beta} u) v}_{\text{SS}} - \frac{1}{2} \int_{\Omega} (\underline{\beta} \cdot \underline{\beta}) uv + \frac{1}{2} \int_{\Omega} (\underline{\beta} \cdot \underline{\beta}) uv$$

In strong form:

$$L = -\mu \Delta u + (\underline{\beta} \cdot \nabla u) + \sigma u = -\mu \Delta u + \sigma u - \frac{1}{2} (\nabla \cdot \underline{\beta}) u + + (\underline{\beta} \cdot \nabla u) + \frac{1}{2} (\nabla \cdot \underline{\beta}) u$$

L_S

L_{SS} .

or

$$L = -\mu \Delta u + \nabla \cdot (\underline{\beta} u) + \sigma u = -\mu \Delta u + \sigma u + \frac{1}{2} (\nabla \cdot \underline{\beta}) u + + \nabla \cdot (\underline{\beta} u) - \frac{1}{2} (\nabla \cdot \underline{\beta}) u$$

L_S

L_{SS}

Notice that $\nabla \cdot (\underline{\beta} u) - \frac{1}{2} (\nabla \cdot \underline{\beta}) u = \frac{1}{2} \nabla \cdot (\underline{\beta} u) + \frac{1}{2} \underline{\beta} \cdot \nabla u$

So also for differential operators (or bilinear forms), we can identify a symmetric and a skew-symmetric part.

Strongly Consistent Stabilization Methods

Let's try to "design" a stabilization method from scratch.

$$\alpha(u_h, v_h) + b_h(u_h, v_h) = \delta(v_h) + g_h(v_h)$$

To maintain the strong consistency, we need that for u , exact solution:

$$(*) \quad b_h(u, v_h) - g_h(v_h) = 0 \quad \forall v_h \in V_h$$

But we have also:

$$(**) \quad \alpha(u, v_h) - \delta(v_h) = 0.$$

So $(*)$ should be related to $(**)$ but how?

Clearly, $(*)$ cannot be $(**) \&$ because this is strongly consistent but adds no stabilization.

(1) → What if we consider the $(**) \&$ in its strong form?

$$\int_{\Omega} (-\mu \Delta u + \beta \cdot \nabla u + \sigma u - f) v \quad (***)$$

This is formally not correct, because the functions are not regular enough!
(In general, u is not H^2).

However, recall that the SD was $\int (\beta \cdot \nabla u) \beta \cdot \nabla v$. This means that if we choose as test function $(\beta \cdot \nabla v)$, we have a potentially stabilizing term.

(2) OK, let's break $(***)$ by elements:

$$\sum_{T_K} \delta_{Ku} \int_{T_K} (-\mu \Delta u_K + \beta \cdot \nabla u_K + \sigma u_K - f) w_K$$

Formally this is correct because on each element the functions are regular.

We are close to the final value.

$$\text{Let's define } b_h(u_h, v_h) - g_h(v_h) = \sum_K \delta_{Ku} \int_{T_K} ((L_S + L_S^S) v_h - f) (L_S v_h + \gamma L_S^S v_h)$$

where δ_{Ku} is a parameter (defined by the user), γ is a parameter with values $\gamma = 0, 1, -1$.

Notice that:

(1) The term is strongly consistent, being proportional to the residual (elementwise, \approx its strong form).

(2) The method includes the SD, since $\sum_k \delta_k \int_{T_k} L u_h \cdot \nabla v_h$ includes exactly the SD terms.

We are legitimated to hope that this is a strongly consistent stabilization.

For $\gamma = 0$ this is called SUPG (Streamline Upwind Petrov-Galerkin)

For $\gamma = 1$ " " " GALS (Galerkin Least Squares)

For $\gamma = -1$ " " " DW (Douglas-Wang).

Remark

For piecewise linear finite elements $\nabla u_h = 0$ on each element, so the actual terms are simpler.

For instance, in the case of $\nabla \cdot \beta = 0$ with linear form we have:

$$a(u_h, v_h) + \sum_k \delta_k \int_{T_k} (\underbrace{\beta \cdot \nabla u_h - \sigma u_h - f,}_{\text{SD}} \underbrace{\beta \cdot \nabla v_h + \gamma \sigma v_h}_{\text{strongly consistent terms}}) = \int_{\Omega} f v_h$$

Analysis of Strongly Consistent Methods

The strong consistency should enable to get results similar to the ones for the regular Galerkin approach. This is true, but we need to introduce the right norms (the strong consistency is only elementwise).

Let's consider the case of GALS ($\gamma = 1$)

We may consider the bilinear form:

$$a_{\text{GALS}}(u_h, v_h) = a(u_h, v_h) + \sum_k \delta_k \left(L u_h, L v_h \right)_{L^2(T_k)}$$

It is possible to prove the following results:

(1) a_{GALS} is continuous in T_k

(2) $a_{\text{GALS}}(u_h, u_h) \geq \mu_0 \| \nabla u_h \|_{L^2}^2 + \gamma_0 \| u_h \|_{L^2}^2 + \sum_k \delta_k \| L u_h \|_{L^2}^2$

$$\mu_0 = \min_{x \in \Omega} \mu(x) \quad \gamma_0 = \min \left((\nabla \cdot \beta) + \frac{1}{2} \sigma \right)$$

Consequently, we can introduce the norm:

$$\|u_\delta\|_{\text{GALS}}^2 = \mu_0 \|\nabla u_\delta\|_{L^2}^2 + \gamma_0 \|u_\delta\|_{L^2}^2 + \sum_k \delta_k \|Lu_\delta\|_{L^2(\Omega_k)}^2$$

The GALS problem is continuous and coercive w.r.t. this norm, so that it is well posed.

Moreover, when $P_k > 1$, for an appropriate choice of δ_k , we can prove:

$$\|u - u_\delta\|_{\text{GALS}} \leq C h^{r+\frac{1}{2}} \|u\|_{H^{r+1}} \quad r = \min(p, s)$$

In this norm, the accuracy depends again on the degree of the FE polynomials.

The parameters δ_k depend on the constant C_0 of the inverse inequality

$$\sum_k h_k^2 \|\Delta u_\delta\|_{L^2(\Omega_k)}^2 \leq C_0 \|\nabla u_\delta\|_{L^2(\Omega)}^2$$

This is called "inverse inequality" because while in general in $H^1(\Omega)$ we have the second derivative dominating the first one with an h of the denominator, in discrete spaces V_h we can prove it.

For SUPG and DWR similar results can be proven (but much more complicated). \Rightarrow See the book Quarteroni-Valli, Chap. 8

REMARK 1

Also in dimension 2+ we can interpret the stabilization methods as Petrov-Galerkin methods (oblique projectors) where

$$W_h = V_h \oplus \text{Bubbles}$$

This interpretation leads to different stabilization methods called "sub-grid viscosity". Find a description in the excellent but advanced book:

Ern, Guermond, Theory and Practice of Finite Elements, Springer

REMARK 2

Reaching Dominated problems give similar oscillations as in 1D also in multiple dimensions.

The Quadrature Formula in 2D on a triangle:

$$\int_{\Delta} f = \frac{1}{3} (f(P_0) + f(P_1) + f(P_2)) |A|$$

where
 $|A|$ = area of the triangle

or in 3D on a tetrahedron

$$\int_{\Delta} f = \frac{1}{4} (f(P_0) + f(P_1) + f(P_2) + f(P_3)) |V|$$

where
 $|V|$ = volume of the T

can be used for the mass matrix with lower finite elements.

(in general, on a simplex in n -dimensions:

$$\int_{\Omega} f \varphi_i = \frac{1}{n+1} \left(\sum_{i=0}^n f(P_i) \right) |V| \quad (\text{where } |V| = \int_{\Omega} 1)$$

Clearly, (linear) finite element functions are orthogonal with respect to the scalar product induced by this quadrature formula, and the resulting matrix is diagonal.

$$(T) \int_{\Omega} \varphi_i \varphi_j = \begin{cases} 0 & i \neq j \\ \neq 0 & i = j \end{cases}$$

This fixes the oscillations.

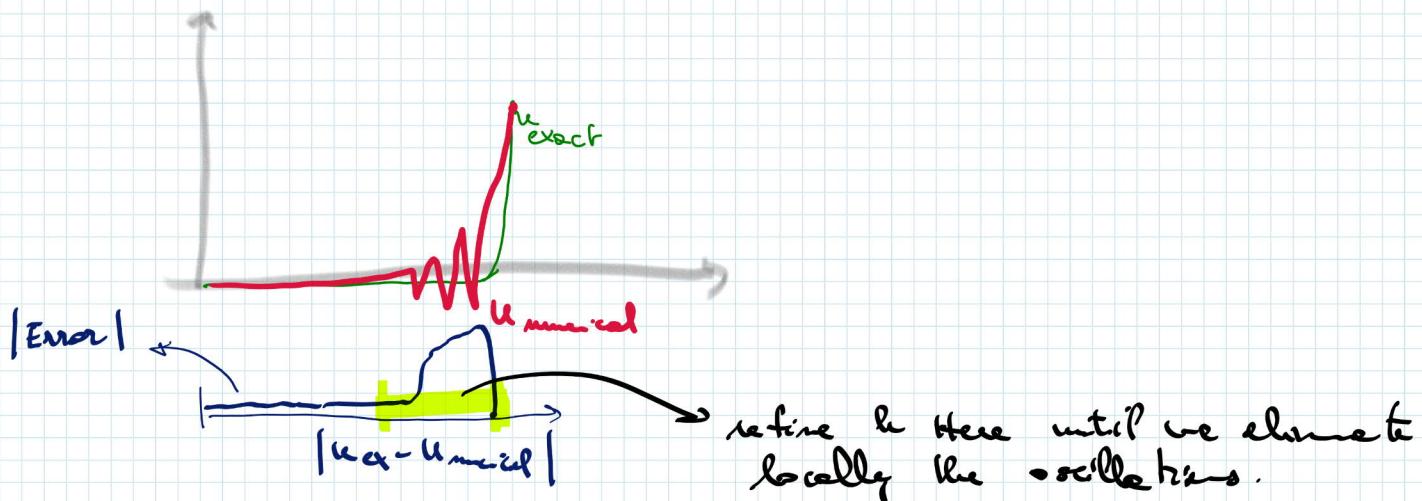
Other (more sophisticated) strategies are needed for higher order elements.

Stabilization and Grid-Adaptivity

An alternative or complementary approach to the Solvechkin / Reaction eliminated problems is to use Grid-Adaptivity.

Since FEM are strongly grid-dependent, we may decide to select the mesh size h in a non-uniform way, so that we can adapt it to the Péclet number.

If β is function of space and we estimate a large local error, we refine h so to make the local $Pe < 1$.



To do this, we need : (i) a good error estimator
(ii) mesh refinement / coarsening methods.

(1) The error estimator theory is huge.

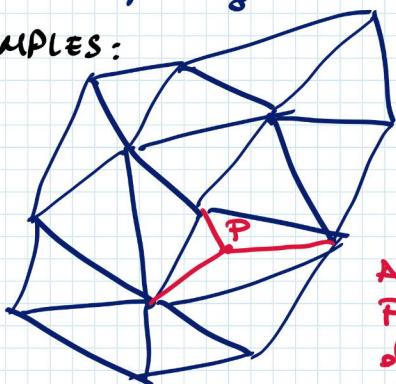
In MATH 516 we consider sometimes some adaptivity and local error estimators are generally based on two solutions with different accuracy. Other approach, like the so-called "goal-oriented" ones are very popular (see e.g. Zienkiewicz-Zhu).

See specific literature (Strikwerda, Rannacher, Brück, etc.)

(2) Refinement and Coarsening are trivial only in 1D.

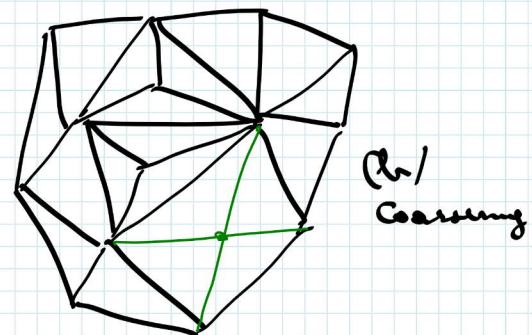
In 2d dimensions, adding or removing a node needs to be done respecting the conformity constraints

EXAMPLES:



(a) Refinement

P
Adding the node
P requires the correct
definition and numbering
of the elements.



(b)
Coarsening

What is the right node
to eliminate for
coarsening?

After removal, how are
organized the elements?

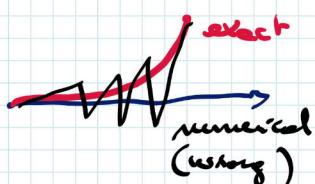
Keep in mind that the great numerical analyst Michel Fortin stated that the mesh should not be considered as given, but as an unknown at the problem.

REMARK

It is well known that for some differential problems we have a principle of maximum or minimum. This means that we know that maximum or minimum are attained only at the boundary.

For the Laplace problem $-\Delta u = 0$, max and min are on $\partial \Omega$.

When we get to the numerical solution, having a similar principle may be useful, because this guarantees that the approximate solution cannot be negative if it is supposed to be positive, or it vanishes on the boundary, for instance.



This requires that in the linear system

$$A \underline{u} = \underline{f}$$

A fulfills some properties that guarantee

The monotonicity of the solution.

This leads to the concept of M-matrix (i.e. $a_{ij} \leq 0$ if $i \neq j$, $\text{Re}(a_{ii}) > 0$)

The construction of a non-oscillatory solution can be characterized by having at the end of the discretization a M-matrix.

Week 11

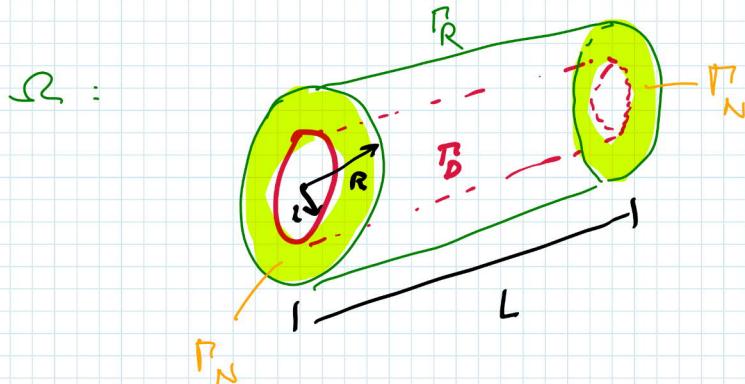
In this week we will learn how to solve the problem

$$-\nabla \cdot (K \nabla u) = 0 \quad \text{in } \Omega$$

$$u(R_D) = u_w$$

$$K \nabla u \cdot \underline{n} (R_N) = 0$$

$$K \nabla u \cdot \underline{n} + \alpha u = \alpha u_{\text{ext}} \quad (T_K)$$



Data

$$R = 0.025 \text{ m}$$

$$r = 0.015 \text{ m}$$

$$L = 0.1 \text{ m}$$

$$u_w = 60^\circ \text{C}$$

$$\alpha = 10 \frac{\text{W}}{\text{m}^2 \text{K}}$$

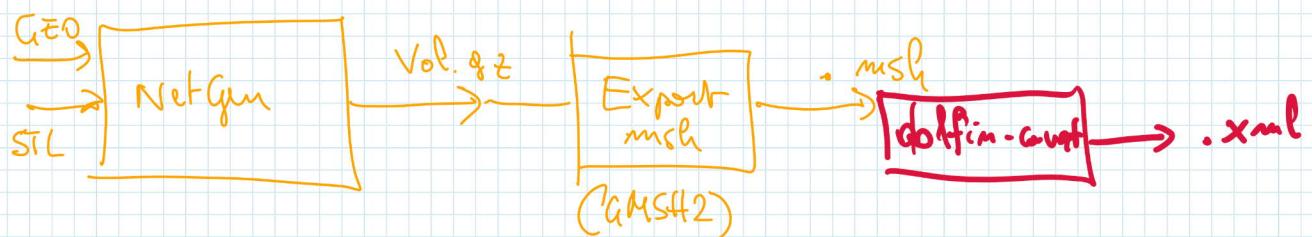
$$K = 0.05 \frac{\text{W}}{\text{m} \text{K}}$$

$$u_{\text{ext}} = 15^\circ \text{C}$$

Preprocessing

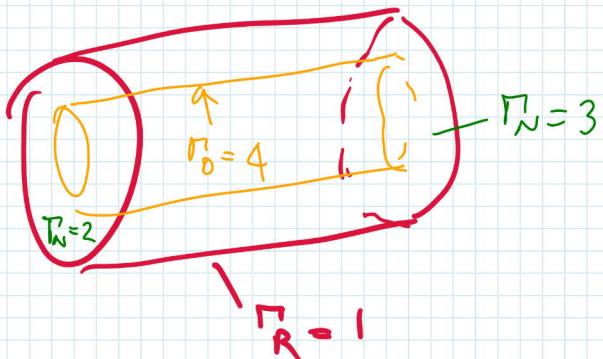
Netgen is a software for meshing with some elementary geometrical subroutines. Generally, an input file is in the STL format, but the constructive geometry format is given in files with extension .geo.

In our case, the geometry is given by two cylinders, internal and external and the internal is subtracted to the external one.



After the mesh is generated with Netgen you need to:

(e) Edit the Boundary Conditions (name "Mesh")



Hint

Use the inverted lexicographic order:

$D < N < R$

↑

Highest index

↑

Lowest Index

In fact, Dirichlet conditions are the LAST ones to be processed in your solver (they are the only ones that affect the pattern of the matrix).

(h) Export to the FENICS format using the GMSH format .msh as an intermediate step.

.geo FILE:

```
#  
## a cylinder  
#  
algebraic3d  
  
# cut cylinder by planes:  
  
solid extcyl = cylinder ( 3, 0, 0; -1, 0, 0; 0.025 )  
and plane (0, 0, 0; -1, 0, 0)  
and plane (0.1, 0, 0; 1, 0, 0);  
  
solid intcyl = cylinder ( 3, 0, 0; -1, 0, 0; 0.015 )  
and plane (-0.1, 0, 0; -1, 0, 0)  
and plane (0.2, 0, 0; 1, 0, 0);  
  
solid insul = extcyl and not intcyl;  
  
tlo insul;
```

Parabolic Problems with FEM

Variational formulation and estimates

In this chapter, we consider problems like:

$$\frac{\partial u}{\partial t} - \mu \Delta u + \beta \cdot \nabla u + \sigma u = f \quad \subset \Omega \times [0, T]$$

+ B.C. like:

$$u(\Gamma_D) = u_0, \quad \mu \nabla u \cdot \underline{n} (\Gamma_{N,D}) = \alpha_N$$

$$\mu \nabla u \cdot \underline{n} + \chi u (\Gamma_{N,R}) = \alpha_R$$

+ I.C.

$$u(x, 0) = u_0(x)$$

$$\text{where } \underline{x} = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

To study this problem, already investigated with the Finite-Difference method with Finite Elements, we need to give a variational formulation.

Notice that in 1D we already gave a variational formulation, without mentioning it. Here, we generalize to multiple dimensions and with the right functional spaces. The Quarteroni's book is only for $\beta=0$, $\sigma=0$, here we carry out a more general analysis.

Let's proceed formally. For the moment, we assume to have only Dirichlet conditions,

$$u(\partial \Omega) = 0$$

So, we postulate $u \in H_0^1(\Omega)$ and we specify the regularity in fine later.

Let $v \in H_0^1(\Omega)$ function, time-independent. With the usual steps, we multiply the equation by v , integrate over Ω .

$$\int_{\Omega} \frac{\partial u}{\partial t} v$$

$$+ \int_{\Omega} \nabla u \cdot \nabla v + \int_{\Omega} \beta \cdot \nabla u v + \int_{\Omega} \sigma u v = \int_{\Omega} f v$$

$$u(\underline{x}, 0) = u_0(\underline{x})$$

if Ω is time-independent:

$$\frac{d}{dt} \int_{\Omega} u v$$

(if Ω is not time-dependent, we need a strong collocated Reynolds theorem:

$$\int_{\Omega} \frac{\partial u}{\partial t} v = \frac{d}{dt} \int_{\Omega} u v - \int_{\partial \Omega} \underline{w} \cdot \underline{n} u v \quad \underline{w} = \frac{\partial \underline{u}}{\partial t}$$

From now on Ω is time-invariant.

So we can write the problem as:

Find $u \in H_0^1(\Omega)$ (time derivative to define) s.t.

$$\frac{d}{dt} \int_{\Omega} u v + a(u, v) = f(v) \quad \forall v \in H_0^1(\Omega) \quad (*)$$

\downarrow
 $\int_{\Omega} f v$
+ I.C.

Lax-Milgram for Parabolic Problems

A sufficient condition for $(*)$ to be well-posed is the continuity of the bilinear form and functionals and:

$$\exists \lambda > 0, \exists \alpha > 0 \text{ s.t. } a(u, u) + \lambda \|u\|_{L^2}^2 \geq \alpha \|u\|_{H^1}^2.$$

This condition is called **WEAK COERCIVITY**, for $\lambda=0$ we get the old coercivity, this is less restrictive.

REMARK Notice that a problem like $\frac{\partial u}{\partial t} - \mu \Delta u - |\sigma| u = f$ in this case is immediately well posed (while with no time-derivative, the Lax-Milgram lemma requires some assumptions on the coefficients)

To complete the picture, we need to specify the functional spaces.

Let us integrate in time the $(*)$:

$$\int_0^T \frac{d}{dt} \int_{\Omega} u v + \boxed{\int_0^T a(u, v)} = \int_0^T f v \rightarrow \begin{aligned} &\text{we need } f \in L^2 \text{ in both time and space} \\ &\text{INTUITION: we need the } \|u\|_{L^2(\cdot, \cdot)} \text{ to be } L^2(0, T) \end{aligned}$$
$$\boxed{\int_{\Omega} u(\xi, T) v - \int_{\Omega} u_0(\xi) v} - \boxed{(u(T), v)_{L^2}} - \boxed{(u_0, v)_{L^2}}$$

INTUITION: we need $\|u\|_{L^2(\cdot, \cdot)}$ to be bounded

We will confirm the intuitions with the estimates on the solution.

For now, let's introduce the spaces:

$$L^2(0, T; H_0^1(\Omega)) = \text{functions with } \|g\|_{H^1(\cdot)} \in L^2(0, T) \text{ (and trace null on } \partial\Omega)$$
$$L^\infty(0, T; L^2(\Omega)) = " g \text{ with } \|g\|_{L^2(\Omega)}(t) \in L^\infty(0, T).$$

$L^2(0, T; L^2(\Omega))$ = functions g s.t. $\|g\|_{L^2(\Omega)} \in L^2(0, T)$.

VARIATIONAL FORMULATION with non-homogeneous Dirichlet conditions.

Find $u \in L^\infty(0, T; L^2(\Omega)) \cap L^2(0, T; H_B^1(\Omega))$ s.t.

given $f \in L^2(0, T; L^2(\Omega))$

$u_0(x) \in L^2(\Omega)$

$u_0(\vec{x}) \in L^2(0, T; H^{1/2}(\partial\Omega))$

$$\frac{d}{dt}(u, v)_{L^2(\Omega)} + a(u, v) = f(v) \quad \forall v \in H_B^1(\Omega)$$

$$u(0) = u_0$$

or

$$(u(t), v)_{L^2(\Omega)} + \int_0^T a(u, v) = \int_0^T f(v) + (u_0, v)_{L^2(\Omega)}$$

A-priori estimates

To get a-priori estimates, we need to go back to the form:

$$\int_{\Omega} \frac{\partial u}{\partial t} v + a(u, v) = \int_{\Omega} f v$$

Now, take $v = u$ (THIS IS TIME-DEPENDENT!):

$$\int_{\Omega} \frac{\partial u}{\partial t} u + a(u, u) = \int_{\Omega} f u$$

$$\text{Now: } \frac{\partial u}{\partial t} u = \frac{1}{2} \frac{\partial (u^2)}{\partial t} \quad \text{so} \quad \int_{\Omega} \frac{\partial u}{\partial t} u = \frac{1}{2} \frac{d}{dt} \|u\|_{L^2(\Omega)}^2$$

We have therefore:

$$\frac{1}{2} \frac{d}{dt} \|u\|_{L^2(\Omega)}^2 + a(u, u) = \int_{\Omega} f u$$

If we integrate in time:

$$\|u\|_{L^2(\Omega)}^2(\tau) + 2 \int_0^\tau a(u, u) = 2 \int_0^\tau \int_{\Omega} f u + \|u_0\|_{L^2(\Omega)}^2$$

Let's recall the Yang Inequality: for α, b given, $\varepsilon > 0$ arbitrary

$$0 \leq \left(\varepsilon\alpha - \frac{b}{2\varepsilon}\right)^2 = \varepsilon^2\alpha^2 + \frac{b^2}{4\varepsilon^2} - \alpha b$$

$$\alpha b \leq \varepsilon^2\alpha^2 + \frac{b^2}{4\varepsilon^2}$$

Then:

$$\int_{\Omega} fu \leq \varepsilon^2 \|u\|_{L^2(\Omega)}^2 + \frac{1}{4\varepsilon^2} \|f\|_{L^2(\Omega)}^2$$

so we can write (under weak coercivity)

$$\frac{1}{2} \frac{d}{dt} \|u\|_{L^2(\Omega)}^2 + \alpha \|u\|_{H^1(\Omega)}^2 - 2\|u\|_{L^2(\Omega)}^2 \leq \varepsilon^2 \|u\|_{L^2(\Omega)}^2 + \frac{1}{4\varepsilon^2} \|f\|_{L^2(\Omega)}^2$$

If we can find ε small enough s.t. $\alpha + \varepsilon^2 < \alpha$, we have:

$$\|u^2\|_{L^2(\Omega)} + 2(\alpha - (\alpha + \varepsilon^2)) \int_0^T \|u\|_{H^1}^2 \leq \frac{1}{4\varepsilon^2} \int_0^T \|f\|_{L^2(\Omega)}^2 + \|u_0\|_{L^2(\Omega)}^2$$

For $\alpha - (\alpha + \varepsilon^2) \geq k > 0$ we immediately get:

$$\begin{aligned} \|u\|_{L^\infty(L^2)} &\leq \text{data} & (T \text{ only replaced with any instant}) \\ \|u\|_{L^2(H^1)} &\leq \text{data} \end{aligned}$$

For a bilinear form coercive ($\lambda = 0$) this is an option i.e. we can choose $\varepsilon^2 = \frac{\alpha}{2}$.

In general, if we cannot say that $\alpha - (\alpha + \varepsilon^2) > 0$, then we can use the Gronwall Lemma:

$$\|u\|_{L^2(\Omega)}^2 + 2\alpha \int_0^t \|u\|_{H^1}^2 \leq \int_0^t \|f\|^2 + \int_0^t (\lambda + 1) \|u\|_{L^2}^2 d\tau + \|u_0\|_{L^2}^2$$

Gronwall lemma:

$$g(t) \leq \alpha + \int_0^t \beta g(\tau) d\tau \Rightarrow g(t) \leq \alpha + \int_0^t \alpha \beta e^{\int_\tau^t \beta} d\tau$$

REMARK

Notice that for u i.e. in $L^2(\Omega)$, the solution for $t > 0 \in L^2(H^1)$, so it is more regular: the diffusion operator gives more regularity to the solution (discontinuities are smoothed).

Semi-Discretization in Space

Let's start using FEM for the derivative in space.

To do this, we postulate for the solution a form similar to the one we use with the separation of variables in 1D:

$$u(x, t) = \sum_{j=1}^{\infty} X_j(x) T_j(t) \quad (*)$$

Here we assume:

$$u_h(x, t) = \sum_{j=1}^N u_j(t) \varphi_j(x) \quad (**)$$

where

$\varphi_j(x)$ are the piecewise polynomial (Lagrange) functions

$u_j(t)$ are functions to find.

The major differences between (*) and (**) are that:

- (i) (*) is a series, (**) is a FINITE sum
- (ii) the X_j are eigenfunctions of the problem, here the $\varphi_j(x)$ are generic polynomial functions with no specific design for the problem at hand.

At this point, we can use (**) into our problem in the form:

$$\frac{d}{dt} (u_h, v_h)_{L^2} + \alpha (u_h, v_h) = f(v_h) \quad \text{where } v_h \text{ is a finite element function in } V_h$$

More precisely, we write:

$$\sum_{j=1}^N \frac{d}{dt} (u_j(t) (\varphi_j, \varphi_i))_{L^2} + \sum_{j=1}^N u_j(t) \alpha (\varphi_j, \varphi_i) = f(\varphi_i) \quad i=1 \dots N$$

that becomes:

$$\sum_j (\varphi_j, \varphi_i)_{L^2} \frac{du_j}{dt} + \sum_j \alpha (\varphi_j, \varphi_i) u_j = f(\varphi_i) \quad i=1 \dots N$$

At the bottom line, this reads

$$M \frac{du}{dt} + Au = b$$

where $u = [u_j]$, M is the mass matrix $[(\varphi_j, \varphi_i)]$,
 A is the stiffness matrix $[\alpha (\varphi_j, \varphi_i)]$
 b is the vector $[f(\varphi_i)]$.

Not surprisingly, we obtain an ordinary differential equation system.

Remark The name "mass matrix" comes from here, because typically the time-derivative comes from "inertial terms" (e.g.: mass \times acceleration when u is a velocity).

Full Discretization

At this point, we can use our Finite difference methods for the time discretization.

However, methods are typically written for $\frac{du}{dt} = Bu + c$
Here we have the mass matrix:

$$M \frac{du}{dt} = -A u + b \quad \text{with } u_0 = \text{approximation of } u_0 \text{ in } V_h \\ (\text{for instance } u_0 = [u_0(x_0)])$$

Two options:

(1) The mass matrix is s.p.g.l., so we can invert it (at least formally, not in practice)

$$\begin{aligned} \frac{du}{dt} &= -M^{-1} A u + M^{-1} b \\ &\approx Bu + c \end{aligned}$$

At this point, we have ≈ 200 of methods: Runge-Kutta, Adams, etc.

For instance, we can use Forward Euler:

$$\left. \begin{array}{l} \rightarrow \text{time step } \Delta t \\ \rightarrow \text{collocation instants } t^n \end{array} \right\} \frac{1}{\Delta t} (u^{n+1} - u^n) = Bu^n + c^n$$

in practice, this reads:

$$Mu^{n+1} = (M - \Delta t A)u^n + \Delta t b^n$$

Backward Euler:

$$(M + \Delta t A)u^{n+1} = \Delta t b^n \quad \left(\text{here, we assume } A \text{ to be time-independent, but this is not necessarily true} \right)$$

Heun: recall that this is the most basic RK method:

$$\frac{1}{\Delta t} (u^{n+1} - u^n) = \frac{1}{2} \left(Bu^n + B(u^n + \Delta t B u^n + \Delta t c^n) \right) + \frac{1}{2} (c^{n+1} + c^n)$$

(Crank-Nicolson + Forward Euler)

This reads

$$Mu^{n+1} = \left(M + \Delta t A + \frac{1}{2} \Delta t^2 A^2 \right) u^n + \frac{\Delta t^2}{2} A M^{-1} b^n + \frac{1}{2} (b^{n+1} + b^n)$$

In practice, we set:

$$\left\{ \begin{array}{l} Mu = (M - \Delta t A)u^n + \Delta t b^n \\ Mu^{n+1} = \frac{1}{2} (A u^n + A \tilde{u}) + \frac{1}{2} (b^{n+1} + b^n) \end{array} \right.$$

Notice that with the Finite Difference method, we have $M = I$, so an explicit method reads

$$\underline{u}^{n+1} = (I - \Delta t A) \underline{u}^n + \Delta t \underline{b}^n \quad (\text{Forward Euler})$$

$$\begin{aligned} \underline{\tilde{u}} &= (I - \Delta t A) \underline{u}^n + \Delta t \underline{b}^n \\ \underline{u}^{n+1} &= \frac{\Delta t}{2} A(\underline{u}^n + \underline{\tilde{u}}) + \frac{\Delta t}{2} (\underline{b}^{n+1} + \underline{b}^n) \end{aligned} \quad (\text{Hem})$$

In this case, we do not have any system to solve. In the FEM case, we have 1+ systems to solve with the matrix M . The mass matrix is s.p.d., so it can be solved easily with CG, however it is an additional cost. The FEM is computationally more expensive.

(ii) Second Option: if we use RE MASS DUMPING (Linear FEM), the problem of explicit methods is solved, as M is diagonal, so the linear systems are trivially solved.

\Rightarrow The mass lumping is useful not only for reaction dominated problems, but also for time-dependent problems.

The Θ -method

As done for FD, we will focus on the Θ -method.

We get to the method by using a quadrature formula:

$$M \int_{t^n}^{t^{n+1}} \frac{du}{dt} = M(u^{n+1} - u^n) = \int_{t^n}^{t^{n+1}} g \approx (\theta g^{n+1} + (1-\theta)g^n) \Delta t$$

We get:

$$(M + \Delta t \theta A) \underline{u}^{n+1} = (M - \Delta t (1-\theta) A) \underline{u}^n + \Delta t (\theta \underline{b}^{n+1} + (1-\theta) \underline{b}^n)$$

Remark

While the weak coercivity is enough?

Let us consider the following problem - weakly coercive:

$$\left(\frac{\partial u}{\partial r}, v \right) + \alpha(u, v) = f(v)$$

Let's change the unknown: $u \rightarrow e^{-\lambda t} u \equiv w$ ($\lambda > 0$)
We have then:

$$e^{-\lambda t} \left(\frac{\partial u}{\partial r}, v \right) + e^{-\lambda t} \alpha(u, v) = e^{-\lambda t} f(v) \equiv g(v)$$

$$\left(\frac{\partial w}{\partial t}, v \right) + \alpha(w, v) + \lambda(w, v) = g(v) \quad \left(\frac{\partial w}{\partial t} = e^{-\lambda t} \frac{\partial u}{\partial t} - \lambda e^{-\lambda t} u \right)$$

$\hat{\alpha}(w, v)$

$$\hat{\alpha}(w, w) = \underbrace{\alpha(w, w)}_{\downarrow} + \lambda \|w\|_{L^2}^2 \quad \left. \begin{array}{l} \geq \alpha \|w\|_{L^2}^2 \\ \alpha \|w\|_H^2 - \lambda \|w\|_{L^2}^2 \end{array} \right\}$$

Formally, a weakly coercive problem can be reformulated into a coercive one (rounding errors here do not matter).

Analysis of the semi-discrete problem

Based on the previous remark, we can focus on a coercive problem (in the traditional sense).

Now, we have the following problem:

Find $u \in L^2(H') \cap L^\infty(L^2)$:

$$\left(\frac{\partial u}{\partial t}, v \right) + \alpha(u, v) = f(v) \quad \forall v \in V$$

and

Find $u_h \in L^2(\mathcal{V}_h) \cap L^\infty(L^2)$:

$$\left(\frac{\partial u_h}{\partial t}, v_h \right) + \alpha(u_h, v_h) = f(v_h) \quad \forall v_h \in \mathcal{V}_h$$

After taking v_h in the first one and subtraction memberwise, we get:

$$\left(\frac{\partial e_h}{\partial t}, v_h \right) + \alpha(e_h, v_h) = 0$$

$$e_h(0) = u_0 - u_{h0} = \eta_0 \quad \forall v_h \in \mathcal{V}_h.$$

where $e_h \equiv u - u_h$

Now we have:

$$\begin{aligned} \alpha \|e_h\|_{H'}^2 &\leq \alpha(e_h, e_h) = \alpha(e_h, u - w_h) + \alpha(e_h, w_h - u_h) = \\ &= \alpha(e_h, u - w_h) - \left(\frac{\partial e_h}{\partial t}, w_h - u_h \right) \end{aligned}$$

At this point, we need to find a bound for the two terms on the right hand side. The first is simpler:

$$|\alpha(e_h, u - w_h)| \leq M \|e_h\| \|u - w_h\| \leq \frac{\alpha}{2} \|e_h\|^2 + \gamma \|u - w_h\|^2$$

The second one requires much more work (see Quarteroni)

At the end, one can prove that

$$\|e_h\|_{L^2(T)}^2 + \int_0^T \|e_h\|_{H'}^2 \leq C \left(\|u\|_{H^s}, \|u\|_{H^s(H^{s+1})} \right) \quad \text{with } s = \min(p, q), \quad u \in H^s(H^{s+1})$$

Stability Analysis of the Θ -method

In this case, we work as done for the FD method:

- discretize in time
- diagonalize the problem
- write the diagonal problem as a set of independent equations

At this point, we find again that the method

$$(M + \Theta \Delta t A) \underline{u}^{n+1} = (M - (1-\Theta) \Delta t A) \underline{u}^n + \Delta t \Theta \underline{b}^n + \Delta t (1-\Theta) \underline{b}^n$$

^{as}

- unconditionally stable for $\Theta \geq \frac{1}{2}$
- conditionally stable for $\Theta < \frac{1}{2}$ with $\Delta t \leq C h^2$

Notice that, in this case, the eigenvalues have to be instead the generalized eigenvalues:

$$A\underline{x} = \lambda M \underline{x} \quad (\text{eigenvalues of } M^{-1} A).$$

Analysis of the Fully Discrete Problem

In this case, we proceed by combining all the previous results for the error $\underline{\epsilon}^n = [(\underline{u}_{\text{fully discrete}} - \underline{u}_{\text{semi-discrete}})(\underline{x}_j, t^n)]$

The proof is technical and long. We need many tricks. For instance, for $\underline{e}, \underline{b}$ two generic vectors:

$$\begin{aligned} (\underline{e} - \underline{a}, \underline{b}) &= (\underline{e} - \underline{a}, \underline{b} - \underline{a}) + (\underline{e} - \underline{a}, \underline{a}) = \\ &= \frac{1}{2} \|\underline{b} - \underline{a}\|^2 + \frac{1}{2} (\underline{e} - \underline{a}, \underline{b} - \underline{a}) + (\underline{e}, \underline{a}) - (\underline{a}, \underline{a}) = \\ &= \frac{1}{2} \|\underline{b} - \underline{a}\|^2 + \frac{1}{2} \|\underline{e}\|^2 - \frac{1}{2} \|\underline{a}\|^2 \end{aligned}$$

$$\text{So that } (\underline{\epsilon}^{n+1} - \hat{\underline{\epsilon}}, \underline{\epsilon}^{n+1}) = \frac{1}{2} \|\hat{\underline{\epsilon}} - \underline{\epsilon}^n\|^2 + \frac{1}{2} \|\underline{\epsilon}^{n+1}\|^2 - \frac{1}{2} \|\underline{\epsilon}^n\|^2$$

and summing over n :

$$\frac{1}{2} \|\hat{\underline{\epsilon}} - \underline{\epsilon}^0\|^2 + \frac{1}{2} \|\underline{\epsilon}^1\|^2 - \frac{1}{2} \|\underline{\epsilon}^0\|^2$$

$$\frac{1}{2} \|\underline{\epsilon}^2 - \underline{\epsilon}^1\|^2 + \frac{1}{2} \|\underline{\epsilon}^3\|^2 - \frac{1}{2} \|\underline{\epsilon}^1\|^2$$

...

$$\sum_{n=1}^{\infty} (\underline{\epsilon}^{n+1} - \hat{\underline{\epsilon}}, \underline{\epsilon}^{n+1}) = \frac{1}{2} \|\hat{\underline{\epsilon}} - \underline{\epsilon}^0\|^2 + \frac{1}{2} \sum_n \|\underline{\epsilon}^{n+1} - \underline{\epsilon}^n\|^2 - \frac{1}{2} \|\underline{\epsilon}^0\|^2$$

Assemblying all the results, we obtain that:

$$\|\underline{e}^n\|_2^2 + \alpha \sum \Delta t \|\underline{e}^n\|_{H^1}^2 \leq C(\text{data}) \left(h^{ex} + \Delta t^{\sigma(\theta)} \right)$$

where $\alpha = \min(\rho, s)$, $\sigma(\theta) = \begin{cases} 1 & \theta \neq \frac{1}{2} \\ 2 & \theta = \frac{1}{2} \end{cases}$.

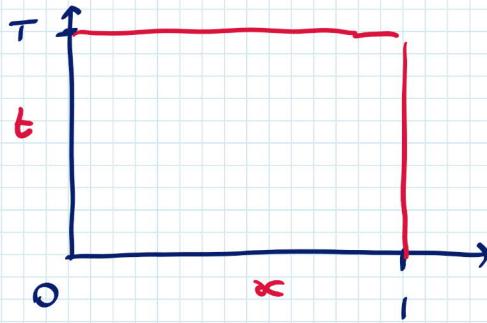
Space-Time Finite Elements

Let us consider the 1D problem:

$$\frac{\partial u}{\partial t} - \mu \frac{\partial^2 u}{\partial x^2} + \beta \frac{\partial u}{\partial x} + \sigma u = f \quad x \in (0, 1) \\ 0 < t \leq T$$

for simplicity.

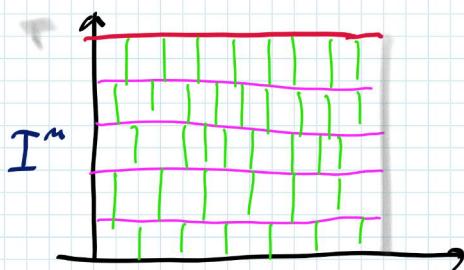
As we said, we could consider space and time at the same way:



However, this leads to a problem where all the time instants are solved simultaneously, which is problematic for simulations with a long time-interval.

We prefer to create a sequential solution in time.

For this reason, we do not mesh the domain with an unstructured mesh but with a sequence of time-slates:



The nodes (green lines) do not need to be collocated in the same places in each slate.

For simplicity, however, we assume they are.

The finite dimensional space ^{in each slate} can be defined as given by functions such that

$$v_n''(x, t) = \left\{ \underbrace{\sum_{k=0}^{\infty} c_k t^k}_{\text{time dependence}} \underbrace{\sum_{j=0}^{\infty} \phi_j \varphi_j(x)}_{\text{space dependence}} \right\}$$

Out of the slate T^n these functions are 0.

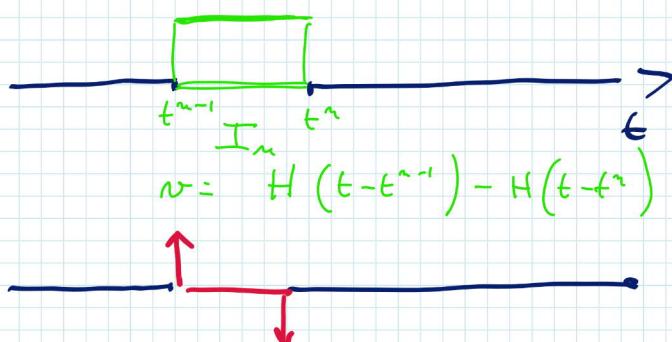
If we now consider the space-time problem:

$$\int_0^T \int_0^1 \frac{\partial u}{\partial t} v + \int_0^T \int_0^1 (-\mu) \frac{\partial^2 u}{\partial x^2} v + \dots = \int_0^T \int_0^1 f dx dt$$

The function $u_h = \sum_{\text{slab}} \sum_{c_i t} c_i t^k \bar{\Sigma} d_j \varphi_j(x)$ is discontinuous in time, so we need to specify the derivative in the distributional sense.

This step requires some technicalities. We cannot write the problem as a sequence of slabs, the would not be independent one to the others, which is wrong.

Notice that if we have a piecewise constant function in time:



the time derivative reads:

$$\frac{du}{dt} = \delta_{t^{n-1}} - \delta_{t^n}$$

For a piecewise linear function: $u = (\alpha + \beta t)\chi$

where χ is the characteristic function of the slab (1 in the slab, 0 out), we have:

$$\frac{du}{dt} = \beta\chi + (\alpha + \beta t)\frac{d\chi}{dt} = \beta\chi + \alpha - (\alpha + \beta\Delta t)$$

↑
the two Dirac delta is real

With these definitions (see Ferriero, Sacco, Vergini) we write the slabwise version of the problem: Chapter 5

$$\int_{I^{n-1}}^1 \int_0^1 \left(\frac{\partial u}{\partial t} v + \mu \frac{\partial u}{\partial x} \frac{\partial v}{\partial x} + \beta \frac{\partial u}{\partial x} v + \tau_{uv} v \right) + \int_0^1 [u^{n-1}] v = \int_0^1 f$$

where

$$[u^{n-1}] = u^{n-1,+} - u^{n-1,-}$$

$\downarrow u^{n-1,+}$
 $\uparrow u^{n-1,-}$ t^{n-1}

It is possible to prove that this formulation is unconditionally stable.

The case with piecewise-constant in time functions.

In this case, in each time slab we assume the solution to be constant in time, so that $u_h(x,t) = \sum u_j^* \varphi_j(x)$

↳ classical Lagrange piecewise polynomials.

So, in this case $u^{n-1,+} = u^n$ and we have $\frac{\partial u}{\partial t} \Big|_{I^n} = 0$.

$$\int_{I^n} \int_0^t \frac{\partial u_e}{\partial t} \varphi_i + \iint_{I^n} \mu \frac{\partial u_e}{\partial x} \frac{\partial \varphi_i}{\partial x} + \beta \frac{\partial u_e}{\partial x} \varphi_i + \sigma u_e \varphi_i = \int_{I^n} \int_0^t f \varphi_i \\ \stackrel{=0}{=} (\text{u.e is constant})$$

↓

$$(u_e^n - u_e^{n-1}, \varphi_i) + \Delta t \Delta (u_e^n, \varphi_i) = \int_{I^n} \int_0^t f \varphi_i$$

If f is 0 or constant in time we have the classical Backward Euler scheme.

The difference, in general, is in the right hand side: the classical BE introduces an additional approximation:

$$\int_{I^n} \int_0^t f \varphi_i \approx \underbrace{(f^n, \varphi_i)}_{\text{Space-Time FEM}} \underbrace{|_{L^2(0,1)}}_{\text{BE}}$$

BE can be regarded as an approximation of Space-Time FEM of order 0 in time.

The case of piecewise linear in time functions.

In this case, we have in each element:

$$u_e(x, t^n) = \alpha_n(x) \psi_0(t) + \beta_n(x) \psi_1(t)$$

where

$$\psi_0 = \frac{t^n - t}{\Delta t} \quad \psi_1 = \frac{t - t^{n-1}}{\Delta t} \quad \left(\frac{d\psi_0}{dt} = -\frac{1}{\Delta t}; \frac{d\psi_1}{dt} = \frac{1}{\Delta t} \right)$$

so that

$$u_e(x, t^{n-1})|_{I^n} = u_e^{n-1, +} = \alpha_n(x)$$

$$u_e(x, t^n)|_{I^n} = u_e^{n, -} = \beta_n(x)$$

So, in the variational formulation we have for the time discretization

$$\left\{ \begin{array}{l} \int_{I^n} \int_0^t \left(u_e^{n-1, +} \frac{d\psi_0}{dt} + u_e^{n, -} \frac{d\psi_1}{dt} \right) \psi_0 v_e + \int_{I^n} \alpha(u_e^{n-1, +} \psi_0 + u_e^{n, -} \psi_1, \psi_0 v_e) \\ + \int_0^t (u_e^{n-1, +} - u_e^{n, -}) v_e = \iint_{I^n} f \psi_0 v_e \\ \int_{I^n} \int_0^t \left(u_e^{n-1, +} \frac{d\psi_0}{dt} + u_e^{n, -} \frac{d\psi_1}{dt} \right) \psi_1 v_e + \int_{I^n} \alpha(u_e^{n-1, +} \psi_0 + u_e^{n, -} \psi_1, \psi_1 v_e) = \int_{I^n} \int_0^t f \psi_1 v_e \end{array} \right.$$

$$\text{Notice that } \int_{I^n} \frac{d\psi_0}{dt} \psi_0 = \int_{t^{n-1}}^{t^n} -\frac{1}{\Delta t^2} (t^n - t) = \left[\frac{1}{2\Delta t^2} (t^n - t)^2 \right]_{t^{n-1}}^{t^n} = -\frac{1}{2}$$

$$\int_{I^n} \frac{d\psi_1}{dt} \psi_0 = \dots = -\frac{1}{2}$$

$$\int_{I^n} \frac{d\psi_0}{dt} \psi_1 = \dots = \frac{1}{2}$$

$$\int_{I^n} \frac{d\psi_1}{dt} \psi_1 = \dots = \frac{1}{2}$$

$$\int_{I^n} \psi_0^2 = \frac{1}{3} \Delta t \quad \int_{I^n} \psi_0 \psi_1 = \frac{1}{6} \Delta t \quad \int_{I^n} \psi_1^2 = \frac{1}{3} \Delta t$$

Proceeding with the discretization, we find the system:

$$\begin{cases} \left(\frac{1}{2} M + \frac{\Delta t}{3} A \right) \underline{v}^{n,-} + \left(\frac{1}{2} M + \frac{\Delta t}{6} A \right) \underline{v}^{n,+} = M \underline{v}^{n-1,-} + \int_{I^n} (\mathbf{f}, \psi_0 \varphi_i)_{\mathbb{C}} \\ \left(-\frac{1}{2} M + \frac{\Delta t}{6} A \right) \underline{v}^{n,-} + \left(\frac{1}{2} M + \frac{\Delta t}{8} A \right) \underline{v}^{n,+} = \int_{I^n} (\mathbf{f}, \psi_1 \varphi_i) \end{cases}$$

At each time step, we need to solve this $(2 \times N) \times (2 \times N)$ linear system.

Accuracy: for linear finite elements:

$$P^0 : \| \mathbf{e} \|_{L^\infty(\mathbb{L}^2)} \leq C (\Delta t \| u \|_{L^\infty(\mathbb{H})} + h^2 \| u \|_{L^\infty(\mathbb{H}^2)})$$

$$P^1 : \| \mathbf{e} \|_{L^\infty(\mathbb{L}^2)} \leq C (\Delta t^2 \| u \|_{L^\infty(\mathbb{H})} + h^2 \| u \|_{L^\infty(\mathbb{H}^2)})$$

Sorry, typo in my book
(as it is written).

Week 13 : Hyperbolic Problems with FEM

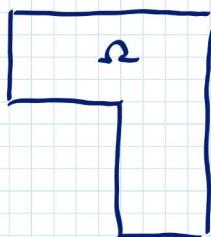
(see the slides)

A short intro to Domain-Decomposition

Domain Decomposition is a technique for the efficient numerical solution of problems :

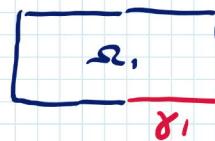
- in complex elements
- with a multiphysics nature
- with parallel processors.

→ surprisingly enough, the first idea of DD is before scientific computing, Schreber in the XIX century :



$$\begin{aligned} -\Delta u = 0 \text{ in } \Omega \\ + \text{B.C.} \end{aligned}$$

⇒



The analytical solution of the Laplace problem is found by separation of variables :

IDEAS

for loop: initial guess for $u(\gamma_1)$

$$\text{solve: } -\Delta u_1^{(u_1)} = 0 + \text{B.C.} + u_1^{(u_1)} = \lambda^{(u_1)}$$

$$\text{solve: } -\Delta u_2^{(u_2)} = 0 + \text{B.C.} + u_2^{(u_2)} = u_1^{(u_1)}(\gamma_2)$$

$$\text{set } \lambda^{(u_1)} = u_2^{(u_2)}(\gamma_1)$$



This is a first example of DD with overlapping.

Notice that if we have an initial guess for $u(\gamma_2)$ too, in █ we can write

$$u_2^{(u_1)}(\gamma_2) = u_1^{(u_1)}(\gamma_2)$$

and the two solution processes are completely independent ⇒ PARALLEL COMPUTING!

One may argue that the conditions used on γ_i ($i=1,2$) do not need to be necessarily Dirichlet.

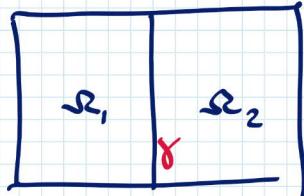
Also, the critical question is: does the iterative process converge?

Chapt. 18 Quasi-Newton

Chapt. 6 Quasi-Newton-Kkt

Chapt. 3 Tomiyama-Selley-V.

(2) Another example is for multi-physics problems:



For instance:

$$-\mu_1 \Delta u_1 = f_1 \quad \text{in } \Omega_1 \\ + \text{B.C.}$$

$$-\mu_2 \Delta u_2 = f_2 \quad \text{in } \Omega_2 \\ + \text{B.C.}$$

+ interface conditions on γ

The interface conditions may reflect physical processes.

In this case we have non-overlapping DD, and it can be used for single physics problem too.

For instance, the Poisson problem:

$$-\Delta u = f \quad \text{in } \Omega \\ + \text{B.C.}$$

is equivalent to

$$-\Delta u_1 = f_1 \quad \text{in } \Omega_1 \\ + \text{B.C.}$$

$$-\Delta u_2 = f_2 \quad \text{in } \Omega_2 \\ + \text{B.C.}_2$$

$$\begin{aligned} u_1(\gamma) &= u_2(\gamma) \\ \nabla u_1 \cdot \mathbf{n}_1(\gamma) &= -\nabla u_2 \cdot \mathbf{n}_2(\gamma) \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{interface conditions.}$$

From here, we have the following non-overlapping DD scheme:

for Loop:

solve $-\Delta u^{(cur)}_1 = f_1$
+ B.C., + $u_1^{(cur)}(\gamma) = \gamma^{(cur)}$

Dirichlet
Neumann

solve $-\Delta u_2^{(cur+1)} = f_2$
+ B.C. + $\nabla u_2 \cdot \mathbf{n}_2(\gamma) = \nabla u_1 \cdot \mathbf{n}_1(\gamma)$

Set $\gamma^{(cur+1)} = u_2^{(cur+1)}(\gamma) \quad (*)$

Again, the critical question is about the convergence.

We will see the convergence analysis.

There are several strategies to improve the convergence.

(1) Replace $(*)$ with $\lambda^{(u)} = \Theta u_2^{(u)}(\gamma) + (1-\Theta) \lambda^{(u)}$
 (Relaxation)

(2) Replace the interface conditions with equivalent Robin-Robin conditions:

$$\nabla u_i^{(u)} + \alpha u_i|_{\gamma} = -\nabla u_2 \cdot n_2 + \alpha u_2|_{\gamma}$$

$$\nabla u_i^{(u)} + \beta u_i|_{\gamma} = -\nabla u_2 \cdot n + \beta u_2|_{\gamma}$$

where α and β are parameters to tune

If they are "optimized", we have "nonoverlapping Schwarz".

As you see, the topic is vast and full of potential:

- preconditioning full problems
- parallel computing and load balancing.

Here, we limit to just some simple analyses in 1D.

1D Analysis

Overlapping Method.

$$\begin{cases} -u_i^{(u)} = 1 \\ u_i(0) = 0 \quad u_i^{(u)}(x_2) = \lambda^{(u)} \end{cases}$$



$$\begin{cases} -u_2^{(u)} = 1 \\ u_2^{(u)}(1) = 0 \quad u_2^{(u)}(x_1) = u_1^{(u)}(x_1) \end{cases}$$

Problem:

$$-u'' = 1$$

$$u(0) = u(1) = 0$$

Test: $\int_{x_1}^{x_2} (u_1^{(u)} - u_2^{(u)})^2 dx$ due then $\lambda^{(u)} = u_2^{(u)}(x_1)$



We know that the exact solution of the problem reads: $\frac{1}{2}x(1-x)$

Let's introduce the errors: $e_1 = u_1 - u$, $e_2 = u_2 - u$

$$\begin{cases} e_1^{(u)} = 0 \\ e_1(0) = 0 \quad e_1^{(u)}(x_2) = e^{(u)} = \lambda^{(u)} - \frac{1}{2}x_2(1-x_2) \end{cases}$$

$$-e_2^{(u)} = 0$$

$$e_2^{(u)}(1) = 0 \quad e_2^{(u)}(x_1) = e_1^{(u)}(x_1)$$

We have:

$$e_1^{(un)} = \epsilon^{(u)} \frac{x}{\alpha_e}$$

↑
linear function

$$e_2^{(un)} = (1-\alpha) \epsilon^{(u)} \frac{x_e}{\alpha_e} \frac{1}{1-\alpha_e}$$

so that $\epsilon^{(un)} = e_2^{(un)}(x_e) = \frac{1-\alpha_e}{1-\alpha} \frac{\alpha_e}{\alpha_e} \epsilon^{(u)}$

To have convergence, i.e. $\epsilon^{(u)} \xrightarrow{k \rightarrow \infty} 0$ we need:

$$\left| \frac{(1-\alpha_e) \frac{\alpha_e}{\alpha_e}}{1-\alpha_e} \right| < 1$$

$$\left| \frac{\frac{1}{\alpha_e} - 1}{\frac{1}{\alpha_e} - 1} \right| < 1 \Rightarrow 1 - \frac{1}{\alpha_e} < \frac{1}{\alpha_e} - 1 < \frac{1}{\alpha_e} - 1$$

↓
 < 0

This is always true for $0 < \alpha_e < \alpha_1 \leq 1$, not for $\alpha_e = \alpha_2$

Also, for either $\alpha_e = 0$ or $\alpha_e = 1$, we have convergence in one step.

This is all expected:

(1) for $\alpha_e = \alpha_1$, we are prescribing two problems with no overlapping, using the same condition twice

(2) for $\alpha_e = 0$ or $\alpha_e = 1$, we are covering the entire domain (take DD).

In general, if $|\alpha_e - \alpha_1|$ gets large (more overlapping), the convergence accelerates:

Study: $\frac{1 - \alpha_e - h}{1 - \alpha_e} \frac{\alpha_e}{\alpha_e h}$ as function of h to see this.

REMARK

This is the so-called "multiplicative version". If we use the parallel or additive version with:

$$u_2^{(un)}(x_e) = u_1^{(K)}(x_e) \xrightarrow{\text{not (un)}} \text{not (un)}$$

we still have the same convergence results, but we need as much as twice iterations (try).

Dirichlet Neumann Method

Let's consider the scheme:

$$\begin{cases} -u_1^{(n+1)} = 1 & x \in [0, \gamma] \\ u_1^{(n+1)}(0) = 0, u_1^{(n+1)}(\gamma) = \lambda^{(n)} \end{cases} \quad \text{when } \gamma \in (0, 1)$$

$$\begin{cases} -u_2^{(n+1)} = 1 \\ u_2^{(n+1)}(1) = 0, \frac{\partial u_2^{(n+1)}}{\partial x} = \frac{\partial u_1^{(n+1)}}{\partial x} \end{cases}$$

$$[\lambda^{(n+1)} = u_2^{(n+1)}(\gamma)]$$

Proceeding with the same notation as before:

$$\begin{cases} -e_1^{(n+1)} = 0 \\ e_1^{(n+1)}(0) = 0, e_1^{(n+1)}(\gamma) = \varepsilon^{(n)} \end{cases} \Rightarrow e_1^{(n+1)} = \frac{\varepsilon^{(n)}}{\gamma} x$$

$$\begin{cases} -e_2^{(n+1)} = 0 \\ e_2^{(n+1)}(1) = 0, e_2^{(n+1)} = \frac{\varepsilon^{(n)}}{\gamma} \end{cases} \Rightarrow e_2^{(n+1)} = -\frac{\varepsilon^{(n)}}{\gamma}(1-x)$$

$$\varepsilon^{(n+1)} = \frac{\varepsilon^{(n)}}{\gamma} (\gamma - 1)$$

The convergence is guaranteed for

$$1 - \gamma < \gamma \Rightarrow \gamma > \frac{1}{2}$$

If we have the Dirichlet domain LARGER than the Neumann elements, we have convergence.

Otherwise, we can use the relaxation:

$$\varepsilon^{(n+1)} = \left(\theta \frac{(\gamma - 1)}{\gamma} + (1 - \theta) \right) \varepsilon^{(n)}$$

$$\downarrow$$

$$\frac{1}{\gamma} \left(\theta \cancel{\gamma} - \theta + \gamma - \cancel{\theta} \gamma \right) \varepsilon^{(n)}$$

$$\left| \frac{\gamma - \theta}{\gamma} \right| < 1 \Rightarrow -\gamma < \gamma - \theta < \gamma$$

↑ trivial

$$\boxed{\theta < 2\gamma}$$

With this choice we converge and the optimal choice is $\theta = \gamma$.

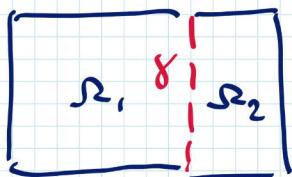
Remark

Computing the flux $\nabla u \cdot \mathbf{n}$ on a Dirichlet boundary.

This topic is somehow overlooked in the literature.

However, this may lead to some useless and inaccurate procedures.
As a matter of fact, computing the flux on a Dirichlet boundary, as required by a DN observer-decomposition can be done with no need of numerical differentiation (that leads to an accuracy degradation).

Let's see:



Dirichlet:

$$-\Delta u_1^{(\text{curl})} = f, \quad + \text{B.C. in } \Sigma_1, \Omega \Sigma_1 \setminus \gamma \\ u_1(\gamma) = g^{(u)}$$

Nemeth:

$$-\Delta u_2^{(\text{curl})} + \text{B.C. in } \Sigma_2, \Omega \Sigma_2 \setminus \gamma \\ \frac{\partial u_2}{\partial \mathbf{n}} = \frac{\partial u_1}{\partial \mathbf{n}} \quad \text{we need this}$$

Let's focus on the Problem 1, "Dirichlet": For simplicity $u(\Omega \Sigma) = 0$

$$-\int_{\Sigma_1} \Delta u_1 v = -\underbrace{\int_{\partial \Sigma} \nabla u_1 \cdot \mathbf{n} v}_{=0 \text{ because } v(\partial \Sigma) = 0} + \int_{\Sigma_1} \nabla u_1 \cdot \mathbf{n} v = \int_{\Sigma_1} f_1 v$$

However, notice that before we apply the Dirichlet boundary conditions, this term is alive.
So, for $v \in H^1(\Sigma_1)$ and not $v \in H_0^1(\Sigma_1)$ we have this term.

$$\text{In } H_0^1(\Sigma_1): \quad \int_{\Sigma_1} \nabla u \cdot \mathbf{n} v = \int_{\Sigma_1} f_1 v \quad v \in H_0^1(\Sigma_1) \\ u \in \mathcal{L} + H_0^1(\Sigma_1)$$

Once we solve this problem, what is:

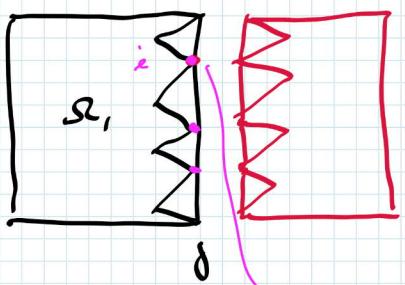
$$\int_{\Sigma_1} \nabla u \cdot \mathbf{n} v - \int_{\Sigma_1} f_1 v = ? \quad \begin{cases} = 0 \text{ for } v \in H_0^1(\Sigma_1) \\ \neq 0 \text{ for } v \in H^1(\Sigma_1) \text{ but } \notin H_0^1(\Sigma_1) \end{cases}$$

The residual is exactly $\int_{\partial \Sigma} \nabla u \cdot \mathbf{n} v$

Now, if we take v vanishing on $\Omega \Sigma_1 \setminus \gamma$, we have

$$\int_{\gamma} \nabla u \cdot \mathbf{n} v$$

This happens with the Lagrange polynomials with nodes on the interface γ .



$$\int_{S_{e_1}} \nabla u \cdot \nabla \varphi_i - \int_{S_{e_1}} f_i \varphi_i = \int_{\delta} \nabla u \cdot n \varphi_i$$

So if we compute

$$\text{Anabc } u - b = \underline{\lambda}$$

\downarrow
No Boundary Conditions

$$\underline{\lambda}(i) \text{ is } \int_{\delta} \nabla u \cdot n \varphi_i$$

If the mesh is conformal (same nodes from left and right) we can use this number as b.c. for S_2 .

Week 14

What I didn't tell you!

Non Linear Problems

When we have a non-linear problem, clearly we need to adjust the procedure with an extra-layer of an iterative root finding step.

Let's see an example:

$$-\Delta u + u^3 = f \quad \text{in } \Omega \\ \text{+ B.C.}$$

Clearly, if we apply the usual procedure, we have to deal with a non-linear algebraic system. We can simplify the problem, by writing

$$u^3 \approx \sum u_i^3 \varphi_i$$

so that we get the non-linear system:

$$K\bar{u} + M\bar{u}^3 = \underline{b} \quad (\bar{u}^3 \text{ in the sense of Matlab } u.^3)$$

At this point, we can call any method for solving nonlinear-algebraic equation. For instance, the Newton method.

$$\underline{F}(u) = 0 \Rightarrow \bar{J} \delta \bar{u}^{(n+1)} = -\underline{F}(u^{(n)}) \quad \bar{J} = \text{jacobian of } \underline{F} \\ \bar{u}^{(n+1)} = u^{(n)} + \delta \bar{u}^{(n+1)} \quad \left(\bar{J}_{ij} = \frac{\partial F_i}{\partial u_j} \right)$$

In the example above:

$$\bar{J}(u) = K + 3M \text{ diag}(u.^2)$$

so we have:

$$\left(K + 3M \text{ diag}(u^{(n)}.^2) \right) (u^{(n+1)} - u^{(n)}) = -Ku^{(n)} - Mu^{(n)}.^3 + \underline{b} \quad \boxed{\text{Linear System}}$$

This needs to be iterated until convergence.

In this approach, we first discretize then linearize. In particular, we approximate

$$\left(\sum u_i \varphi_i \right)^3 \approx \sum u_i^3 \varphi_i$$

This last step is optional.

There is another approach: LINEARIZE then DISCRETIZE.

$-\Delta u + u^3 = f$ can be written as $-\Delta u + u^{*2}u = f$

when $u^* \approx u$. For instance :

- give $u^{(0)}$

- loop k : $-\Delta u^{(k+1)} + (u^{(k)})^2 u^{(k+1)} = f$ until convergence (if possible)

- we discretize this problem :

$$K \underline{u}^{(k+1)} + M(u^{(k)})^2 \underline{u}^{(k+1)} = \underline{f}$$

where $M(u^{(k)})_{ij} = \int \underline{u}^{(k)} \varphi_j \varphi_i$
given

As usual, the problem here is to guarantee the convergence. A relaxation step may help:

$$K \tilde{u} + M(u^{(k)}) \tilde{u} = \underline{f}$$

$$u^{(k+1)} = \theta \tilde{u} + (1-\theta) u^{(k)} \quad \theta \in (0,1]$$

This is called Picard iteration, and it is in general:

- easy
- robust (converges for a large neighbourhood of the solution)
- slow

A Newton method can be used with the ^{Re} Lim-Then-Disce approach, but it requires the differentiation of the operator.

To find the so-called TANGENT PROBLEM we differentiate:

$$\lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} (-\Delta(u + \varepsilon v) + (u + \varepsilon v)^3 + \Delta u - u^3) =$$

$$= \lim_{\varepsilon \rightarrow 0} \left(-\varepsilon \Delta v + u^3 + 3\varepsilon u^2 v + 3\varepsilon^2 u v^2 + \varepsilon^3 v - u^3 \right) =$$

$$> \lim_{\varepsilon \rightarrow 0} (-\Delta v + 3u^2 v + 3\varepsilon u v^2 + \varepsilon^3 v) = -\Delta v + 3u^2 v$$

so the tangent problem reads:

$$-\Delta \delta u + 3u^2 \delta u = f + \Delta u - u^3$$

Iteratively, one has to solve:

given $\underline{u}^{(0)}$;

loop:

$$-\Delta \underline{u}^{(u+1)} + 3\underline{u}^{(u)}^2 \underline{u}^{(u+1)} - 3\underline{u}^{(u)}^3 =$$
$$= f - \underline{u}^{(u)}^3$$

↓,

$$-\Delta \underline{u}^{(u+1)} + 3\underline{u}^{(u)}^2 \underline{u}^{(u+1)} = f + 2\underline{u}^{(u)}^3$$

CHECK: if $\underline{u}^{(u+1)} = \underline{u}^{(u)} = \bar{u}$ then

$$-\Delta \bar{u} + \bar{u} = f \Rightarrow \text{Yes!!!}$$

The Newton method is more delicate, it requires a good initial guess to converge (better than Picard), but it is faster.

A reasonable approach is

$$\underline{u}^{(0)} \rightarrow \text{Picard for a few iterations} \rightarrow \underline{u}_{\text{Picard}} = \underline{u}_{\text{Newton}}^{(0)} \rightarrow \text{Newton.}$$

To summarize:

DISCRETIZE - THEN - LINEARIZE \Rightarrow Nonlinear Algebraic Solver

LINEARIZE - THEN - DISCRETIZE \Rightarrow sequence of nonlinear problems.

The linearization can be done in other ways, for instance fixed point iterations.

$$-\Delta u + \sin(u) u = f$$

could lead to:

$$-\Delta u^{(u+1)} + \sin(u^{(u)}) u^{(u+1)} = f \quad (\text{Picard})$$

$$-\Delta u^{(u+1)} = f - \sin(u^{(u)}) u^{(u)} \quad (\text{Probably not convergent or slow})$$

$$-\Delta u^{(u+1)} + \cos(u^{(u)}) u^{(u)} (u^{(u+1)} - u^{(u)}) + \sin(u^{(u)}) (u^{(u+1)} - u^{(u)}) = f - \sin(u^{(u)}) u^{(u)} \quad (\text{Newton})$$

$$\begin{aligned} & \downarrow \\ -\Delta u^{(u+1)} + \cos(u^{(u)}) u^{(u)} u^{(u+1)} + \sin(u^{(u)}) u^{(u+1)} &= \\ &= f + \cos(u^{(u)}) u^{(u)} u^{(u+1)^2} \end{aligned}$$

Personally, I never Linearize - then - Discretize.

EXAMPLE

Viscous Shearless Burgers Equation (similar to Navier-Stokes)

$$(\underline{u} \cdot \nabla) \underline{u} - \mu \Delta \underline{u} = \underline{f} \quad \underline{u} = \begin{bmatrix} u_0 \\ u_1 \\ u_2 \end{bmatrix}$$

memberwise:

$$(\Sigma_i) u_i \frac{\partial u_j}{\partial x_i} - \mu \sum_i \frac{\partial^2 u_j}{\partial x_i^2} = f_j$$

understood in the Einstein notation

→ $u_0 \frac{\partial u_0}{\partial x_0} + u_1 \frac{\partial u_0}{\partial x_1} + u_2 \frac{\partial u_0}{\partial x_2} - \mu \frac{\partial^2 u_0}{\partial x_0^2} - \mu \frac{\partial^2 u_0}{\partial x_1^2} - \mu \frac{\partial^2 u_0}{\partial x_2^2} = f_0$

Picard: $(\underline{u}^{(u)}, \nabla \underline{u}^{(u+1)}) - \mu \Delta \underline{u}^{(u+1)} = \underline{f}$

Newton: TANGENT PROBLEM:

$$\lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} \left((\underline{u} + \varepsilon \underline{v}) \cdot \nabla (\underline{u} + \varepsilon \underline{v}) - \mu \Delta (\underline{u} + \varepsilon \underline{v}) - \underline{u} \cdot \nabla \underline{u} + \mu \Delta \underline{u} \right) =$$

$$= (\underline{v} \cdot \nabla) \underline{u} + (\underline{u} \cdot \nabla) \underline{v} - \mu \Delta \underline{v}$$

In the Newton iteration, we have ($\delta \underline{u} = \underline{u}^{(u+1)} - \underline{u}^{(u)} = \underline{v}$)

$$(\delta \underline{u} \cdot \nabla) \underline{u}^{(u)} + (\underline{u}^{(u)} \cdot \nabla) \delta \underline{u} - \mu \Delta \delta \underline{u} = \underline{f} - (\underline{u}^{(u)} \cdot \nabla \underline{u}^{(u)} + \mu \Delta \underline{u}^{(u)})$$

$$(\underline{u}^{(u+1)} \cdot \nabla) \underline{u}^{(u+1)} + (\underline{u}^{(u)} \cdot \nabla) \underline{u}^{(u+1)} - \mu \Delta \underline{u}^{(u+1)} = \underline{f} + (\underline{u}^{(u)} \cdot \nabla) \underline{u}^{(u)}$$

REMARK

The FENICS tutorial suggests that there is a DISCRETIZE-THEN-NEWTON approach encoded.

However, Linearize-Then-Discretize approaches are mostly implemented.

The case of Time-Dependent Problems.

In the case of time-dependent problems, non-linearity may introduce some outstanding computational costs.

As a matter of fact, we may have 3 nested loops: !:

$$\left\{ \begin{array}{l} \text{for time} \\ \text{for nonlinear equations} \\ \text{for linear solvers (if iterative)} \end{array} \right.$$

However notice that this cost applies only for implicit time-solvency.

In fact, if we have:

$$\frac{\partial u}{\partial t} - \mu \Delta u + u^3 = f$$

an explicit scheme needs:

$$u^{n+1} = u^n + \Delta t K \tilde{A}^{-1} f(u^n) + \Delta t f^n$$

The implicit one (linearize then discrete)

$$u^{n+1} - \Delta t \mu \Delta u^{n+1, u^n} + \Delta t 3(u^{n+1, u^n})^2 u^{n+1, u^n} = \Delta t f^n + \Delta t 2(u^{n+1, u^n})^3$$

\Rightarrow linear system:

$$\left[I + \Delta t K + \Delta t 3M(u^{n+1, u^n}) \right] u^{n+1, u^n} = \Delta t b + \Delta t 2C(u^{n+1, u^n})$$

An explicit scheme, however, may have poor stability properties.

A possible TRADE-OFF is a SEMI-IMPLICIT SCHEME, where we use the time advancing to linearize the problem:

$$\frac{1}{\Delta t}(u^{n+1} - u^n) - \mu \Delta u^{n+1} + (u^n)^2 u^{n+1} = f^{n+1}$$

RATIONAL: $u^n = u^{n+1} + O(\Delta t)$

$O(\Delta t)$ is also
the discretization
in time error

With a second order scheme we could write:

$$u^{n+1} - u^n \quad \text{where } u^n = 2u^n - u^{n-1}$$

Notice that, in this way, the stability is conditional, but as long as the linearization concerns terms of order O (reactive) or 1 (convective) the constraint $\Delta t < Ch^2$ is avoided.

EXAMPLE FOR BURGERS:

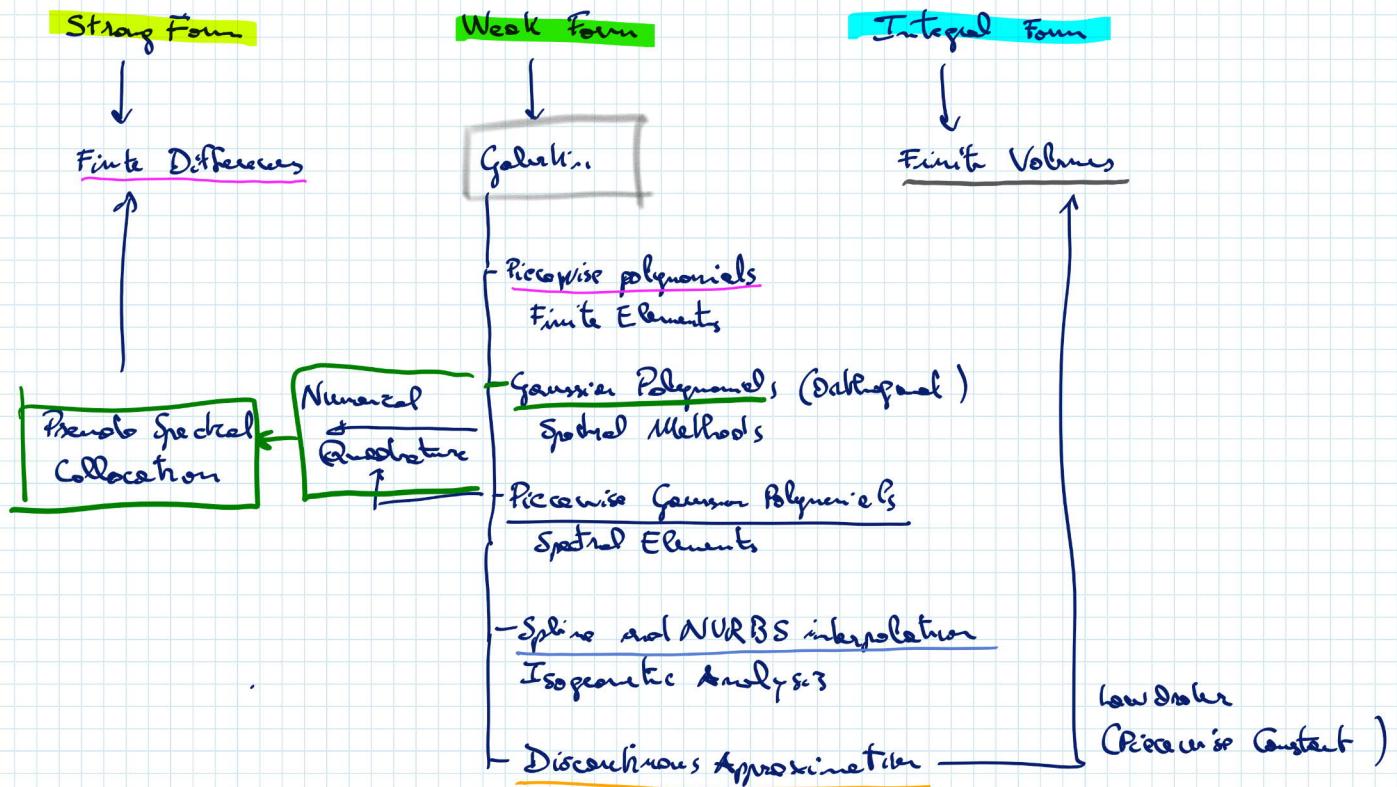
$$\frac{\partial u}{\partial t} - \mu \Delta u + (\underline{u} \cdot \nabla) \underline{u} = f$$

$$\frac{1}{\Delta t}(u^{n+1} - u^n) - \mu \Delta u^{n+1} + (\underline{u}^n \cdot \nabla) \underline{u}^{n+1} = f$$

\underline{u}^n extrapolation in time: \underline{u}^n
(of the same order of the time discretization).

A quick look to other methods

Differential Problems:



Two words on Spectral Methods

When we use Gaussian interpolation as the backbone, we have some advantages and some drawbacks.

Spectral Methods are extremely accurate : the Gaussian interpolation has an exponential accuracy ($N = \text{degree of the chosen polynomials}$)

$$\|u - u_N\|_{H^1} \leq C \|u\|_{H^s} N^{-s} \quad \text{for } u \in H^s$$

$$(s = +\infty \leq C \|u\|_{H^s} e^{-s})$$

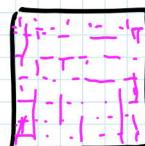
so the only limiting factor is the regularity of the solution

The identification of the nodes for the Gaussian interpolation is easy in 1D and simple geometries. In a complex geometry, it is not trivial (we should map a complex geometry into a simple one)

\Rightarrow this led to the development of Spectral Elements.

Quadrilaterals in Spectral Methods are popular as they are obtained as the tensor-product of the 1D case

For a square
+ - + - + - + -



Spectral basis functions are L^2 -orthogonal (Sobolev, Legendre) \Rightarrow

the mass matrix is diagonal ! $(\varphi_i, \varphi_j) = 0$ for $i \neq j$!

There is no way that for polynomials of order $N \rightarrow +\infty$ we do exact integration \Rightarrow All the integrals are approximated by suitable quadrature formulas.

This means that the method is actually a pseudo-spectral and the analysis tool is the STRANG-LEMMA (not Cea).

If we approximate:

$$(\nabla \varphi_i, \nabla \varphi_j) \approx \sum_{k=1}^{n_{qu}} w_k \nabla \varphi_i \cdot \nabla \varphi_j(x_k, y_k, z_k)$$

$(n_{qu} = \text{number of quadrature nodes})$

it is possible to reinterpret the entire method as a collocation method (Finite Difference) applied to the strong form.

ASIDE NOTE

Finite Elements have been recently extended to polytopes by Brezzi and co-workers (Virtual Element Method or Mimetic Finite Differences).

Saddle Point Problems

When facing elliptic (and also parabolic) problems, we look for the solution of a free minimization:

$$\min_{u \in H_0^1} J = \frac{1}{2} \int_{\Omega} |\nabla u|^2 - \int_{\Omega} f u \Rightarrow -\Delta u = f \text{ in } \Omega \\ u = 0 \text{ on } \partial \Omega$$

This leads to a coercivity-based formulation:

$$a(u, v) = (f, v) \quad \forall v \in H_0^1 \quad \bullet$$

with $a(u, u) \geq \alpha \|u\|^2 \quad \blacksquare$

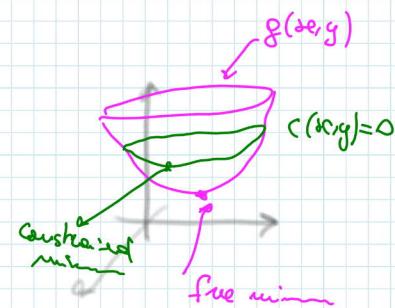
The coercivity plays the role of convexity in minimization:

$$\min_{x \in \mathbb{R}} g \quad (\Leftarrow) \quad g'(x) = 0 \quad \bullet \\ \text{with } g'' > 0 \quad \blacksquare$$

What happens if the minimization is not free?

In calculus we may have:

$$\min g(x, y) \\ \text{with the constraint } c(x, y) = 0$$



In the specific case, one could find

$$c(x, y) = 0 \Rightarrow y = y(x) \quad \boxed{\text{Diri theorem}}$$

Then

$$g(x, y) \Big|_{c(x, y) = 0} = g(x, y(x)) = \tilde{g}(x)$$

and minimize $\tilde{g}(x)$.

To pursue a general strategy, one can introduce the Lagrange multiplier approach:

$$\min g(x, y) \text{ with constraint } c(x, y) = 0$$



$$\max_{\lambda} \min_{x, y} \mathcal{L}(x, y; \lambda) = g(x, y) - \lambda c(x, y) \quad (\text{free})$$

In fact, this leads to:

$$\frac{\partial \mathcal{L}}{\partial \lambda} : c(x, y) = 0 \quad (\text{constraint})$$

$$\frac{\partial \mathcal{L}}{\partial (x, y)} : \frac{\partial \mathcal{L}}{\partial (x, y)} - \lambda \frac{\partial c}{\partial (x, y)} = 0$$

What about problems involving functionals?

AN IMPORTANT EXAMPLE

$$\tilde{J}(\underline{u}) = \frac{1}{2} \int_{\Omega} |\nabla \underline{u}|^2 - \underline{f} \cdot \underline{u} \quad \underline{u} \text{ vector}$$

$$\text{with the constraint } \nabla \cdot \underline{u} = 0$$

$$\Rightarrow \mathcal{L} = \int_{\Omega} \nabla \underline{u} \cdot \nabla \underline{v} - \int_{\Omega} p \nabla \cdot \underline{u} = 0$$

$$\text{for } \underline{u} \in \underline{\mathcal{S}}(\Omega) = \underline{\mathcal{Q}}$$

$$\begin{cases} -\Delta \underline{u} + \nabla p = 0 \\ \nabla \cdot \underline{u} = 0 \end{cases}$$

STOKES PROBLEM
(incompressible fluids)

Now, the coercivity is not enough to be well posedness.
We need something more.

FORMAL DISCUSSION

Find $(\underline{u}, p) \in (H_0^1(\Omega), L^2(\Omega))$ s.t.

$$\begin{cases} \int_{\Omega} \nabla \underline{u} \cdot \nabla \underline{v} - \int_{\Omega} p \nabla \cdot \underline{v} = 0 & \forall \underline{v} \in H_0^1(\Omega) \\ \int_{\Omega} q \nabla \cdot \underline{u} = 0 & \forall q \in L^2(\Omega) \end{cases}$$

$$\text{Find } (\underline{u}, p) \in \mathcal{V} \times \mathcal{Q}: \mathcal{A}(\underline{u}, \underline{v}) + \mathcal{B}(p, q) = 0 \quad \forall \underline{v} \in \mathcal{V}$$

$$\mathcal{B}(\underline{u}, q) = 0 \quad \forall q \in \mathcal{Q}$$

Sufficient conditions:

$$\exists \alpha > 0 \text{ s.t. } \forall \underline{u} \in \mathcal{V} \quad \mathcal{A}(\underline{u}, \underline{u}) \geq \alpha \|\underline{u}\|_{\mathcal{V}}^2$$

inf-sup condition

$$\left[\forall p \in \mathcal{Q}, \exists \underline{v} \in \mathcal{V} \text{ and } \beta > 0 \text{ s.t. } \mathcal{B}(p, \underline{v}) \geq \beta \|p\|_{\mathcal{Q}} \|\underline{v}\|_{\mathcal{V}} \right]$$

These conditions are true for Stokes.

But... but... what happens when we perform a Galerkin approach.

Simply coercive case:

$$\text{coercivity in } V \Rightarrow \text{coercivity in } V_h$$

$$a(u, u) \geq \alpha \|u\|_V^2 \Rightarrow a(u_h, u_h) \geq \alpha \|u_h\|_V^2 \in V_h \subset V$$

In the saddle point case:

$$a(u, u) \geq \alpha \|u\|_V^2 \Rightarrow a(u_h, u_h) \geq \alpha \|u_h\|_V^2 \in V_h \subset V$$

but:

$$\forall p \in L^2, \exists \underline{w} \in T, \beta > 0 \text{ s.t.}$$

$$b(v, p) \geq \beta \|p\|_0 \|v\|_V$$



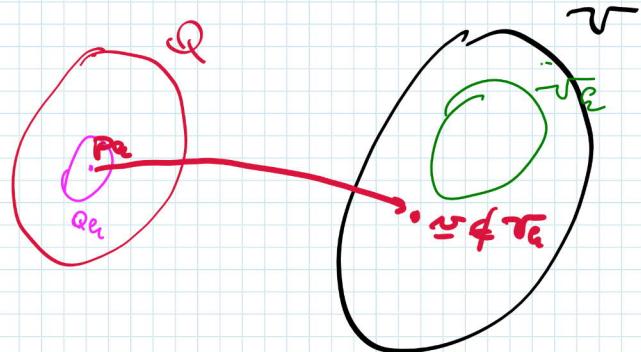
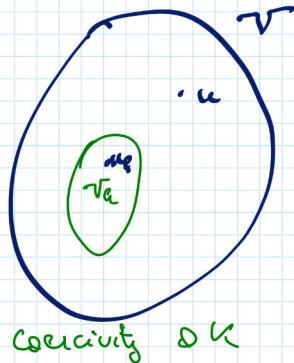
$$\forall p_h \in Q_h, \exists \underline{w}_h \in V_h, \beta > 0 \text{ s.t.}$$

$$b(\underline{w}_h, p_h) \geq \beta \|p_h\|_0 \|v\|_V$$

even when $Q_h \subset Q$ and $T_h \subset T$

This requires some extra effort.

In diagrams:



We need to enforce that V_h contains all the corresponding elements of Q_h for the inf-sup condition.

This excludes some pairs of finite elements.

For instance, for Stokes:

$\underline{u}_h \in P^1, p_h \in P^1$ is NOT inf-sup compatible

$\underline{u}_h \in P^2, p_h \in P^1$ is inf-sup compatible.

