

# Numerical Methods for Differential Equations

Ferran Arqué

2018

Notes for the course *Numerical Methods for Differential Equations* taught by Jaume Soler and Esther Sala at Facultat de Matemàtiques i Estadística of Universitat Politècnica de Catalunya.

## Contents

### Ordinary Differential Equations

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Ordinary Differential Equations. Basic concepts</b> | <b>4</b>  |
| 1.1      | Introduction and some notation                         | 4         |
| 1.2      | Euler's method   | 4         |
| 1.2.1    | Formulation  | 4         |
| 1.2.2    | Local truncation error                                 | 5         |
| 1.2.3    | Convergence  | 5         |
| 1.3      | Improved Euler's method                                | 6         |
| 1.3.1    | Formulation  | 6         |
| 1.3.2    | Local truncation error                                 | 7         |
| 1.4      | Final remarks  | 8         |
| <b>2</b> | <b>Runge-Kutta and Linear Multistep Methods</b>        | <b>9</b>  |
| 2.1      | General Runge-Kutta methods                            | 9         |
| 2.1.1    | Embedded Runge-Kutta methods                           | 10        |
| 2.2      | Linear multistep methods                               | 11        |
| 2.2.1    | Generalities   | 12        |
| 2.2.2    | Richardson's extrapolation                             | 13        |
| 2.2.3    | Convergence of a linear multistep method               | 14        |
| <b>3</b> | <b>Stiff Problems</b>                                  | <b>18</b> |

### Partial Differential Equations

|          |   |           |
|----------|---|-----------|
| <b>4</b> | <b>Partial Differential Equations. Generalities on their solution</b> | <b>25</b> |
| 4.1      | Finite differences  | 25        |
| 4.1.1    | Numerical derivatives   | 26        |
| 4.1.2    | Forward Time Centered Space method (FTCS)                             | 28        |

|          |   |           |
|----------|---|-----------|
| <b>5</b> | <b>Numerical Solution of PDEs with the Finite Elements Method . . . . .</b> | <b>30</b> |
| 5.1      | The Finite Elements Method (FEM) . . . . .                                  | 30        |
| 5.1.1    | From the strong form to the weak form . . . . .                             | 30        |
| 5.1.2    | Discretisation . . . . .  | 31        |
| 5.2      | FEM analysis . . . . .  | 38        |

This work is licensed under a [Creative Commons](https://creativecommons.org/licenses/by-nc-sa/4.0/) “Attribution-NonCommercial-ShareAlike 4.0 International” license.



# ORDINARY DIFFERENTIAL EQUATIONS

# 1 Ordinary Differential Equations. Basic concepts

## 1.1 Introduction and some notation

Given  $y' = f(x, y)$ , where  $\begin{cases} y(x) \in \mathbb{R}^n \\ f : \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}^n \end{cases}$

**Definition.** We denote by  $y(x)$  the exact solution of the ODE system above.

**Definition.**  $y_k$  is the approximation of  $y(x_k)$  (after  $k$  steps).

**Objective.** We want to approximate  $y(x)$  within a given interval  $[x_0, x_n]$ .

We know  $\begin{cases} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{cases}$       We'd like to know  $\begin{cases} y(x_0) \\ y(x_1) \\ y(x_2) \\ \vdots \\ y(x_n) \end{cases}$       We find  $\begin{cases} y_0 \\ y_1 \\ y_2 \\ \vdots \\ y_n \end{cases}$  (given by a method)

**Definition.**  $\|y(x_n) - y_n\|$  is the **global error**.

**Definition.** We define the **local truncation error** as the error caused by one iteration, i.e.

$$LTE = \|y(x_k) - y_k\| \quad (\text{assuming the } \textit{localizing assumption}: y_{k-1} = y(x_{k-1}))$$

**Definition.** A method is of **order**  $p$  if the LTE is of the form

$$LTE = Ch^{p+1} + \mathcal{O}(h^{p+2})$$

## 1.2 Euler's method

### 1.2.1 Formulation

Consider the differential equation system

$$y'(x) = f(y(x)), \quad y(x_0) = y_0 \quad (\text{with } f : \mathbb{R}^n \rightarrow \mathbb{R}^n) \quad (1.1)$$

We'll also assume, for later results, that  $f$  satisfies a Lipschitz condition

$$\|f(y) - f(z)\| \leq L \|y - z\|$$

Now, we want to find an approximation of (1.1) on an interval  $[x_0, x_N]$ , so we'll create a set of  $N + 1$  evenly spaced points between  $x_0$  and  $x_N$ , where  $N$  is the number of steps and  $h = \frac{x_N - x_0}{N}$  the stepsize.

To do that, we can use Euler's method

$$y_{n+1} = y_n + hf(y_n)$$

### 1.2.2 Local truncation error

We want to find  $LTE = \|y(x_{n+1}) - y_{n+1}\|$ . Let's expand it

$$\begin{aligned} y(x_{n+1}) - y_{n+1} &= y(x_n + h) - y_{n+1} = y(x_n) + hy'(x_n) + \mathcal{O}(h^2) - y_n - hf(y_n) \stackrel{\text{Loc. ass.}}{=} \\ &= y(x_n) + hy'(x_n) + \mathcal{O}(h^2) - y(x_n) - hf(y(x_n)) \\ &= hy'(x_n) + \mathcal{O}(h^2) - hy'(x_n) = \mathcal{O}(h^2) \end{aligned}$$

So  $LTE = \mathcal{O}(h^2)$ , and the method has order  $p = 1$ .

### 1.2.3 Convergence

For the method to be convergent, it has to verify

$$\|y(x_N) - y_N\| \xrightarrow{h \rightarrow 0} 0$$

Similar to what we did with the LTE, we'll find a bound for the error in one step, but this time without using the localizing assumption:

$$\begin{aligned} \|y(x_{n+1}) - y_{n+1}\| &= \|y(x_n + h) - y_n - hf(y_n)\| \stackrel{\text{Taylor}}{=} \|y(x_n) + h \underbrace{y'(x_n)}_{f(y(x_n))} + \tilde{c}h^2 - y_n - hf(y_n)\| \\ &\leq \|y(x_n) - y_n\| + h \underbrace{\|f(y(x_n)) - f(y_n)\|}_{L \wedge \frac{1}{L\|y(x_n) - y_n\|}} + ch^2 \leq (1 + hL) \|y(x_n) - y_n\| + ch^2 \quad \forall n \end{aligned}$$

Thus

$$\begin{aligned} \|y(x_1) - y_1\| &\leq (1 + hL) \|y(x_0) - y_0\| + c_1 h^2 \\ \|y(x_2) - y_2\| &\leq (1 + hL) \|y(x_1) - y_1\| + c_2 h^2 \leq (1 + hL)^2 \|y(x_0) - y_0\| + (1 + hL)c_1 h^2 + c_2 h^2 \\ \|y(x_3) - y_3\| &\leq (1 + hL)^3 \|y(x_0) - y_0\| + (1 + hL)^2 c_1 h^2 + (1 + hL)c_2 h^2 + c_3 h^2 \\ &\vdots \end{aligned}$$

Let  $C = \max_i c_i$

Then

$$\begin{aligned} \|y(x_N) - y_N\| &\leq (1 + hL)^N \|y(x_0) - y_0\| + \underbrace{\left(1 + (1 + hL) + (1 + hL)^2 + \dots + (1 + hL)^{N-1}\right)}_{\substack{\text{geom} \rightarrow \| \\ \frac{(1+hL)^N - 1}{hL}}} Ch^2 \\ \implies \|y(x_N) - y_N\| &\leq (1 + hL)^N \underbrace{\|y(x_0) - y_0\|}_{\substack{(1.1) \| \\ 0}} + \frac{(1 + hL)^N - 1}{hL} \cdot Ch^2 \leq Ah \rightarrow 0 \end{aligned}$$

**Remark.** In the last inequality we used  $(1 + hL)^N = e^{N \log(1+hL)} \leq_{\substack{\log(1+x) \\ \leq x+cx^2}} e^{N(hL+ch^2L^2)} =$   
 $= e^{L(x_N - x_0)} \cdot e^{\tilde{\alpha}(x_N - x_0)h} = \tilde{A} \cdot e^{\alpha h} \rightarrow \tilde{A} \cdot 1 = \tilde{A} \quad (Nh = x_N - x_0)$

## 1.3 Improved Euler's method

### 1.3.1 Formulation

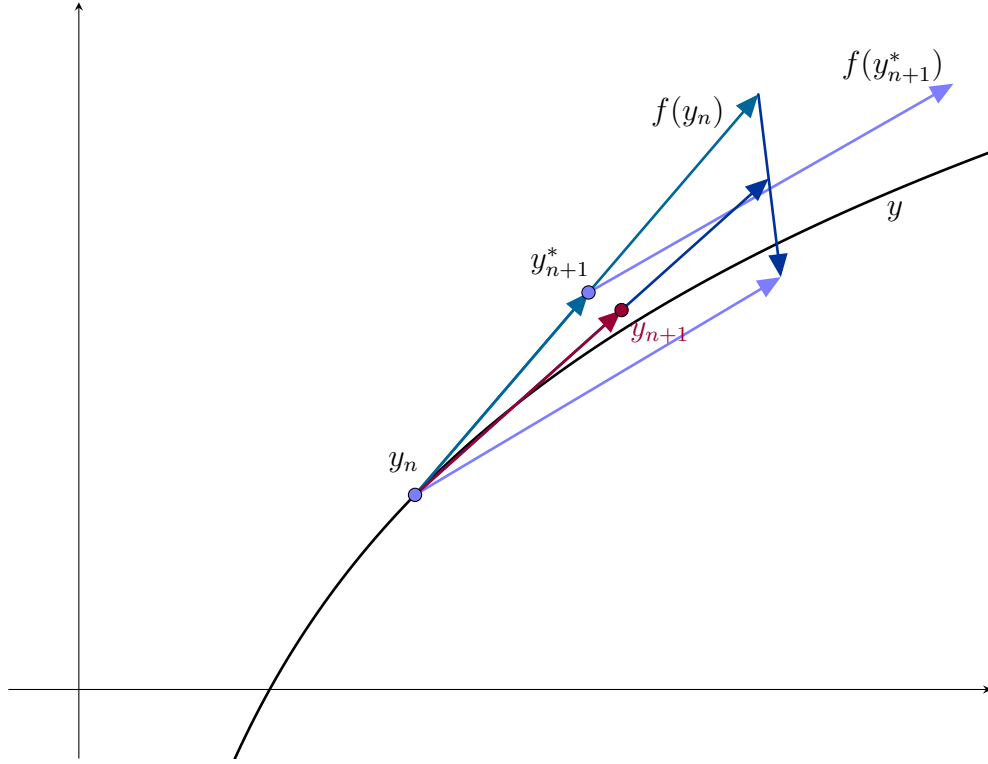


Figure 1.1: One step of the enhanced Euler's method

Given  $y'(x) = f(y(x))$ ,  $y : \mathbb{R} \rightarrow \mathbb{R}$ ,  $f : \mathbb{R} \rightarrow \mathbb{R}$ , we'll go through the steps to deduce the improved Euler's method (also called Heun's method) with the help of the scheme in Figure 1.1.

The auxiliary point  $y_{n+1}^*$  can be found doing one step of the standard Euler's method, so

$$y_{n+1}^* = y_n + h \cdot f(y_n)$$

To get the point  $y_{n+1}$  we compute the average vector of  $f(y_{n+1}^*)$  and  $f(y_n)$ , and with this new vector, we can apply again a step of Euler's method, ending up with our method

$$y_{n+1} = y_n + \frac{h}{2} \cdot (f(y_{n+1}^*) + f(y_n))$$

### 1.3.2 Local truncation error

The Local Truncation Error is given by:

$$LTE = \|method - exact\ solution\|$$

Then

$$LTE = \|y_{n+1} - y(x_{n+1})\| = \left\| y_n + \frac{h}{2}f(y_n) + \frac{h}{2}f(y_{n+1}^*) - \underbrace{y(x_{n+1})}_{y(x_n+h)} \right\| = (*)$$

Applying Taylor on  $f(y_{n+1}^*) = f(y_n + hf(y_n))$ , we have

$$f(y_n + hf(y_n)) = f(y_n) + hf(y_n)f'(y_n) + \mathcal{O}(h^2)$$

and on  $y(x_n + h)$  we have

$$y(x_n + h) = y(x_n) + h \cdot y'(x_n) + \frac{h^2}{2} \cdot y''(x_n) + \mathcal{O}(h^3)$$

(In this case, we expand to second order for later simplifications)

$$(*) = \left\| y_n + \frac{h}{2}f(y_n) + \frac{h}{2}(f(y_n) + hf(y_n)f'(y_n) + \mathcal{O}(h^2)) - \left( y(x_n) + h \cdot y'(x_n) + \frac{h^2}{2} \cdot y''(x_n) + \mathcal{O}(h^3) \right) \right\| \quad (1)$$

Now, given  $y'(x) = f(y(x))$ , we have

$$\begin{aligned} y''(x) &= f'(y(x)) \cdot y'(x) \\ &= f'(y(x)) \cdot f(y(x)) \end{aligned}$$

and we can rewrite the following expression as:

$$\frac{h^2}{2}f(y(x))f'(y(x)) = \frac{h^2}{2}y''(x)$$

With that and the localising assumption ( $y_n = y(x_n)$ ), we can simplify most of the terms in (1) and we end up with

$$\left\| \cancel{y_n} + \cancel{\frac{h}{2}f(y_n)} + \cancel{\frac{h}{2}f(y_n)} + \cancel{\frac{h^2}{2}f(y_n)f'(y_n)} + \mathcal{O}(h^3) - \left( \cancel{y(x_n)} + \cancel{h \cdot y'(x_n)} + \cancel{\frac{h^2}{2} \cdot y''(x_n)} + \mathcal{O}(h^3) \right) \right\| = \mathcal{O}(h^3)$$

So  $LTE = \mathcal{O}(h^3)$

**Remark.** Of course, this method also works for  $y : \mathbb{R} \rightarrow \mathbb{R}^n$ ,  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$

## 1.4 Final remarks

- There's an enhanced Euler's method of order 2 similar to the previous one:

$$y_{n+1} = y_n + hf\left(\frac{y_{n+1}^* + y_n}{2}\right)$$

- If we have a method of order  $\geq 2$ , and we want the value of  $y(x^*)$ , with  $x^*$  off the mesh ( $x^* \neq k \cdot h$ ), we have some options:

-Step back and take a step with the right  $h$ .

-Interpolate with the right order.

-Use a continuous Runge-Kutta method.



## 2 Runge-Kutta and Linear Multistep Methods

### 2.1 General Runge-Kutta methods

Runge-Kutta methods are a family of iterative methods, which include the previously seen Euler's Method. Let's define an RK method with  $s$  stages:

Given  $x \in \mathbb{R}$ ,  $y \in \mathbb{R}^n$ ,  $y' = f(x, y)$

**Definition.** The **Butcher Tableau** is

|          |          |                   |          |          |
|----------|----------|-------------------|----------|----------|
| $c_1$    | $a_{11}$ | $a_{12}$          | $\dots$  | $a_{1s}$ |
| $c_2$    | $a_{21}$ | $a_{22}$          |          | $\vdots$ |
| $\vdots$ | $\vdots$ |                   | $\ddots$ | $\vdots$ |
| $c_s$    | $a_{s1}$ | $\dots\dots\dots$ |          | $a_{ss}$ |
|          | $b_1$    | $\dots\dots\dots$ |          | $b_s$    |

With these coefficients given by the table, we can now define our Runge-Kutta method:

$$\left. \begin{aligned} k_1 &= f\left(x_n + c_1 h, y_n + h(a_{11}k_1 + a_{12}k_2 + \dots + a_{1s}k_s)\right) \\ k_2 &= f\left(x_n + c_2 h, y_n + h(a_{21}k_1 + a_{22}k_2 + \dots + a_{2s}k_s)\right) \\ &\vdots \\ k_s &= f\left(x_n + c_s h, y_n + h(a_{s1}k_1 + a_{s2}k_2 + \dots + a_{ss}k_s)\right) \end{aligned} \right\} \begin{array}{l} \text{System of equations} \\ \text{with unknowns } k_1, k_2, \dots, k_s \end{array}$$

$$\boxed{y_{n+1} = y_n + h(b_1k_1 + b_2k_2 + \dots + b_sk_s)}$$

Cases:

- (1) Explicit ( $a_{ij} = 0$  for  $j \geq i$  and  $c_1 = 0$ )
- (2) Semi-implicit ( $a_{ij} = 0$  for  $j > i$ )
- (3) Implicit

**Remark.** (2) and (3) are used for stiff problems.

**Theorem 2.1.1.** An explicit  $s$ -stage Runge-Kutta method can't have order  $> s$ .

**Theorem 2.1.2.** There is no explicit 5-stage RK of order 5.

**Theorem 2.1.3.**

Let

$A = \text{order } p \text{ for } y' = f(y), f : \mathbb{R}^m \rightarrow \mathbb{R}^m, m > 1$

$B = \text{order } p \text{ for } y' = f(x, y), f : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$

$C = \text{order } p \text{ for } y' = f(y), f : \mathbb{R} \rightarrow \mathbb{R}$

Then

For  $1 \leq p \leq 3$ ,  $A \iff B \iff C$

For  $p = 4$ ,  $A \iff B \implies C$  but  $C \not\Rightarrow B$

For  $p \geq 5$ ,  $A \implies B \implies C$  but  $C \not\Rightarrow B$ ,  $B \not\Rightarrow A$

### 2.1.1 Embedded Runge-Kutta methods

Embedded methods combine two methods of order  $p$  and  $p + 1$  to obtain an estimate of the LTE. This error estimate is useful for adaptive stepsize methods.

In this case, the *Butcher tableau* is

|          |          |                   |             |               |             |
|----------|----------|-------------------|-------------|---------------|-------------|
| 0        |          |                   |             |               |             |
| $c_2$    | $a_{21}$ |                   |             |               |             |
| $c_3$    | $a_{31}$ | $a_{32}$          |             |               |             |
| $\vdots$ | $\vdots$ |                   | $\ddots$    |               |             |
| $c_s$    | $a_{s1}$ | $\dots\dots\dots$ | $a_{s,s-1}$ |               |             |
|          |          |                   |             | $b_1$         | $b_s$       |
|          |          |                   |             | $\bar{b}_1$   | $\bar{b}_s$ |
|          |          |                   |             | order $p$     |             |
|          |          |                   |             | order $p + 1$ |             |

This kind of method has the advantage that it usually requires less steps even though there are more evaluations for each step.

Matlab's current `ode45` solver uses an embedded RK method with orders 4 and 5.

## 2.2 Linear multistep methods

**Definition.** A linear k-step method is a method of the form

$$\sum_{j=0}^k \alpha_j y_{n+j} = h \sum_{j=0}^k \beta_j f_{n+j}$$

with  $\alpha_j, \beta_j \in \mathbb{R}$ ,  $\alpha_k = 1$  and  $\alpha_0 \neq 0$  or  $\beta_0 \neq 0$ .

Notation:  $f_m = f(x_m, y_m)$ ,  $x_m = x_0 + mh$

**Remark.** If  $\beta_k = 0$ , then the method is explicit.

**Remark.** With these methods, adapting the step size can be complicated. It may even change the order.

### Example 2.2.1

Our initial value problem is

$$\begin{cases} y' = f(x, y) \\ y_0 = y(x_0) \end{cases}$$

And we want to solve it using the following multistep method:

$$y_{n+2} - y_n = \frac{h}{2}(f_{n+2} + 4f_{n+1} + f_n)$$

How do we compute  $y_1$ ?

Usually we first use a one-step method like Euler's to compute the initial terms needed, and then we apply the multistep method.

### Example 2.2.2

Let's see how we compute the order of  $y_{n+2} - y_n = h(\beta_1 f_{n+1} + \beta_0 f_n)$ :

$$\text{LTE} = \|y(x_{n+2}) - y_{n+2}\|$$

$$\begin{aligned} y(x_n + 2h) - y_{n+2} &= \cancel{y(x_n)} + 2hy'(x_n) + \frac{1}{2}y''(x_n)4h^2 + \frac{1}{3!}y'''(x_n)8h^3 + \mathcal{O}(h^4) - \\ &\quad - \left( \cancel{y(x_n)} + h\beta_1 f(y_{n+1}) + h\beta_0 f(y_n) \right) \\ &= 2hy'(x_n) + \frac{1}{2}y''(x_n)4h^2 + \frac{1}{3!}y'''(x_n)8h^3 + \mathcal{O}(h^4) - \\ &\quad - \left( h\underbrace{\beta_1 y'(x_n + h)}_{y'(x_n) + hy''(x_n) + \dots} + h\beta_0 y'(x_n) \right) \end{aligned}$$

Grouping by powers of  $h$  we get:

$$h \longrightarrow 2y'(x_n) - \beta_1 y'(x_n) - \beta_0 y'(x_n) = 0 \implies \beta_0 + \beta_1 = 2$$

$$h^2 \longrightarrow 2y''(x_n) - \beta_1 y''(x_n) = 0 \implies b_1 = 2 \implies \beta_0 = 0$$

$$\text{So LTE} = \mathcal{O}(h^3)$$

### 2.2.1 Generalities

- **Adams–Bashforth methods:** These are explicit methods with coefficients  $\alpha_{k-1} = -1$  and  $\alpha_{k-2} = \alpha_{k-3} = \dots = \alpha_0 = 0$ . The  $\beta_j$ 's are chosen such that the method has order  $k$  ( $j = 0, \dots, k$ ).

For example, the Adams–Bashforth method with  $k = 2$  is

$$y_{n+2} = y_{n+1} + h \left( \frac{3}{2} f_{n+1} - \frac{1}{2} f_n \right)$$

- **Adams–Moulton methods:** These have the same coefficients as the Adams–Bashforth methods, but they are implicit methods. Also, for a  $k$ -step Adams–Moulton method, if  $\beta_k \neq 0$ , it can reach order  $k + 1$ , while an A-B method has only order  $k$ .

The Adams–Moulton method with  $k = 2$  is

$$y_{n+2} = y_{n+1} + h \left( \frac{5}{12} f_{n+2} + \frac{2}{3} f_{n+1} - \frac{1}{12} f_n \right)$$

- **Backward differentiation formulas (BDF):** These are implicit methods with  $b_{k-1} = \dots = b_0 = 0$  and the rest of the coefficients are chosen such that the method has order  $k$  (only exist for  $k < 7$ ). BDF methods are useful for stiff problems.
- **Predictor–Corrector method:** A Predictor–Corrector method uses an explicit method for the predictor step (P) and an implicit method for the corrector step (C). They usually have the same order.

$$y_n \xrightarrow{P} y_{n+1}^* \xrightarrow{C} y_{n+1}$$

They have the advantage that the LTE is easy to estimate ( $|predictor - corrector|$ )  
The improved Euler's method (1.3) is an example of a P-C method.

### 2.2.2 Richardson's extrapolation

Richardson's extrapolation is used to improve the rate of convergence of a method.

Suppose we have a one-step explicit method. Then

- Compute one step with stepsize  $h$ :  $y_n \rightarrow y_{n+1}$
- Compute two steps with stepsize  $\frac{h}{2}$ :  $y_n \rightarrow y_{n+\frac{1}{2}} \rightarrow \hat{y}_{n+1}$

Then

$$\text{LTE} \simeq \|y_{n+1} - \hat{y}_{n+1}\|$$

If the method is of order  $p$ , we have

$$y(x_{n+1}) - y_{n+1} = Ch^{p+1} + \mathcal{O}(h^{p+2}) \quad (2.1)$$

$$y(x_{n+1}) - \hat{y}_{n+1} = 2C \left(\frac{h}{2}\right)^{p+1} + \mathcal{O}(h^{p+2}) \quad (2.2)$$

And subtracting (2.1) - (2.2), we end up with

$$\hat{y}_{n+1} - y_{n+1} = C \left( h^{p+1} - 2 \frac{h^{p+1}}{2^{p+1}} \right) = Ch^{p+1} \underbrace{\left( 1 - \frac{1}{2^p} \right)}_{\frac{2^p - 1}{2^p}} \simeq \text{LTE}$$

and

$$C = \frac{\|\hat{y}_{n+1} - y_{n+1}\|}{h^{p+1}} \cdot \frac{2^p}{2^p - 1}$$

Now, while applying our method, we want to keep the LTE below a tolerance  $TOL$ . If the LTE is larger than  $TOL$ , we want to reduce the stepsize, and if it's significantly smaller, we don't want to waste computational time, so we should increase the stepsize. So an automatic adjustment of  $h$  can be computed as follows:

Let  $h^*$  be a new stepsize such that the LTE is equal to the tolerance ( $C \cdot (h^*)^{p+1} = TOL$ ).

Then

$$(h^*)^{p+1} = \frac{TOL}{C} = \frac{TOL}{\|\hat{y}_{n+1} - y_{n+1}\|} \cdot \frac{2^p - 1}{2^p} h^{p+1}$$

And the new stepsize is given by

$$h^* = 0.9h \left( \frac{TOL}{\|\hat{y}_{n+1} - y_{n+1}\|} \cdot \frac{2^p - 1}{2^p} \right)^{\frac{1}{p+1}}$$

### 2.2.3 Convergence of a linear multistep method

**Theorem 2.2.1.** Given a general linear multistep method

$$y_{n+k} + \alpha_{k-1}y_{n+k-1} + \dots + \alpha_0y_n = h(\beta_k f_{n+k} + \dots + \beta_0 f_n)$$

and the polynomials

$$\begin{aligned}\rho(z) &= z^k + \alpha_{k-1}z^{k-1} + \dots + \alpha_0 \\ \sigma(z) &= \beta_k z^k + \beta_{k-1}z^{k-1} + \dots + \beta_0\end{aligned}$$

Then, a necessary and sufficient condition for this method to be convergent is:

1. It is of order 1 at least.
2. The roots of the first stability polynomial ( $\rho(z)$ ) are on the unit ball, and if they are on the boundary, they are simple (as roots).

**Remark.** 1 is always a root of  $\rho(z)$ .

Let's see an example of divergence using a linear multistep method:

#### Example 2.2.3

*Given the method*

$$y_{n+2} + a_1y_{n+1} + a_0y_n = h(b_1f_{n+1} + b_0f_n)$$

1) Find  $a_0, a_1, b_0, b_1$  so that the method above has the highest possible order.

2) Try it on

$$\begin{cases} y' = -y \\ y(0) = 1 \end{cases} \quad (y_0 = 1, y_1 = e^{-h})$$

*and prove the method diverges.*

1) We want  $y(x_n + 2h) - y_{n+2}$

We assume  $y_{n+1} = y(x_n + h), y_n = y(x_n)$  (localizing assumption).

$$\begin{aligned}y(x_n + 2h) - y_{n+2} &= y(x_n + 2h) - \left[ -a_1y_{n+1} - a_0y_n + h(b_1f(y_{n+1}) + b_0f(y_n)) \right] \Big|_{\text{loc.as.}} \\ &= y(x_n + 2h) - \left[ -a_1y(x_n + h) - a_0y(x_n) + hb_1 \underbrace{f(y(x_n + h))}_{y'(x_n+h)} + hb_0 \underbrace{f(y(x_n))}_{y'(x_n)} \right]\end{aligned}$$

As usual, we expand in powers of  $h$ . We'll expand to order 3

$$\begin{aligned}
& y(x_n) + 2hy'(x_n) + \frac{4h^2}{2}y''(x_n) + \frac{8h^3}{6}y'''(x_n) + o(h^4) - \\
& - \left[ -a_1 \left( y(x_n) + hy'(x_n) + \frac{h^2}{2}y''(x_n) + \frac{h^3}{6}y'''(x_n) + o(h^4) \right) - \right. \\
& \quad - a_0y(x_n) \\
& \quad + hb_1 \left( y'(x_n) + hy''(x_n) + \frac{h^2}{2}y'''(x_n) + o(h^3) \right) + \\
& \quad \left. + hb_0y'(x_n) \right]
\end{aligned}$$

Let's group by powers of  $h$  and assume the right conditions to obtain the highest possible order:

$$h^0 \longrightarrow y(x_n) + a_1y(x_n) + a_0y(x_n) = 0$$

$$h^1 \longrightarrow 2hy'(x_n) + a_1hy'(x_n) - hb_1y'(x_n) - hb_0y'(x_n) = 0$$

$$h^2 \longrightarrow 2h^2y''(x_n) + a_1\frac{1}{2}h^2y''(x_n) - b_1h^2y''(x_n) = 0$$

$$h^3 \longrightarrow \frac{8h^3}{6}y'''(x_n) + a_1\frac{h^3}{6}y'''(x_n) - b_1h \left( \frac{h^2}{2}y'''(x_n) \right) = 0$$

With that, we get the system of equations

$$\begin{cases}
1 + a_1 + a_0 = 0 \\
2 + a_1 - b_1 - b_0 = 0 \\
2 + \frac{a_1}{2} - b_1 = 0 \\
\frac{8}{6} + \frac{a_1}{6} - \frac{b_1}{2} = 0
\end{cases}$$

And we end up with

$$a_0 = -5, \quad a_1 = 4, \quad b_0 = 2, \quad b_1 = 4$$

2) Our method is

$$y_{n+2} + 4y_{n+1} - 5y_n = h(4f_{n+1} + 2f_n)$$

and with

$$\begin{cases} y' = -y \\ y(0) = 1, y(h) = e^{-h} \end{cases} \quad (y(x) = e^{-x})$$

we have

$$y_{n+2} + 4y_{n+1} - 5y_n = h(-4y_{n+1} - 2y_n)$$

We'll find a solution of the form

$$y_n = c_1(\quad)^n + c_2(\quad)^n$$

and we'll see that it diverges.

$$\begin{aligned} \lambda^2 + 4\lambda - 5 + 4h\lambda + 2h &= 0 \\ \lambda^2 + (4(1+h))\lambda + (2h-5) &= 0 \end{aligned}$$

$$\lambda = \frac{-4(1+h) \pm \sqrt{4^2(1+h)^2 - 4(2h-5)}}{2}$$

Let's expand the discriminant

$$\begin{aligned} \sqrt{4^2(1+2h+h^2) - 8h + 20} &= \sqrt{36 + 24h + 16h^2} = 6\sqrt{1 + \frac{4}{6}h + \frac{4^2}{6^2}h^2} \stackrel{\text{Taylor}}{=} \\ &= 6\left(1 + \frac{1}{2}\left(\frac{4}{6}h + \frac{4^2}{6^2}h^2\right) + o(h^2)\right) = \\ &= 6\left(1 + \frac{1}{3}h + o(h^2)\right) \end{aligned}$$

So

$$\lambda = \frac{-4 - 4h \pm (6 + 2h + o(h^2))}{2} = \begin{cases} 1 - h + o(h^2) \\ -5 - 3h + o(h^2) \end{cases}$$

$$\implies y_n = c_1(1 - h + o(h^2))^n + c_2(-5 - 3h + o(h^2))^n$$



Let's find  $c_1$  and  $c_2$  imposing the initial conditions

$$\begin{aligned}
& \begin{cases} 1 = c_1 + c_2 \implies c_1 = 1 - c_2 \\ e^{-h} = c_1(1 - h + o(h^2)) + c_2(-5 - 3h + o(h^2)) \end{cases} \\
& \implies e^{-h} = (1 - h + o(h^2)) + c_2 \underbrace{(-5 - 3h + o(h^2) - 1 + h + o(h^2))}_{-6 - 2h + o(h^2)} \\
& \implies c_2 = \frac{e^{-h} - 1 + h + o(h^2)}{-6 - 2h + o(h^2)} \stackrel{\text{Taylor } e^{-h}}{=} \frac{1 - h + o(h^2) - 1 + h + o(h^2)}{-6 - 2h + o(h^2)} = \frac{o(h^2)}{-6 - 2h + o(h^2)} \\
& \implies c_1 = 1 + \frac{1}{6 + 2h + o(h^2)} = \frac{7 + 2h + o(h^2)}{6 + 2h + o(h^2)}
\end{aligned}$$

So

$$y_n = \frac{7 + 2h + o(h^2)}{6 + 2h + o(h^2)} (1 - h + o(h^2))^n + \frac{o(h^2)}{-6 - 2h + o(h^2)} (-5 - 3h + o(h^2))^n$$

and the term  $(-5)^n$  will cause the solution to diverge.

### 3 Stiff Problems

In some ODEs, the step size taken by an adaptive method is forced to be unreasonably small even in regions where the solution curve is smooth. In these cases, it takes a large amount of steps to go through a short time interval.

These types of equations are called **stiff ODEs**.

**Example 3.0.1** (Van der Pol equation)

Given the Van der Pol equation

$$\ddot{x} - \mu(1 - x^2)\dot{x} + x = 0 \quad (\text{with } x(0) = 2)$$

the larger the constant  $\mu$ , the stiffer is the problem.

Trying to solve it using an explicit adaptive stepsize method like Matlab's `ode45` yields



Figure 3.1: Van der Pol equation solution with `ode45` ( $\mu = 10$ )

With 873 steps needed, and a minimum stepsize of  $2.5119 \cdot 10^{-5}$ .

Now, using an implicit method like Matlab's `ode15s`, we have

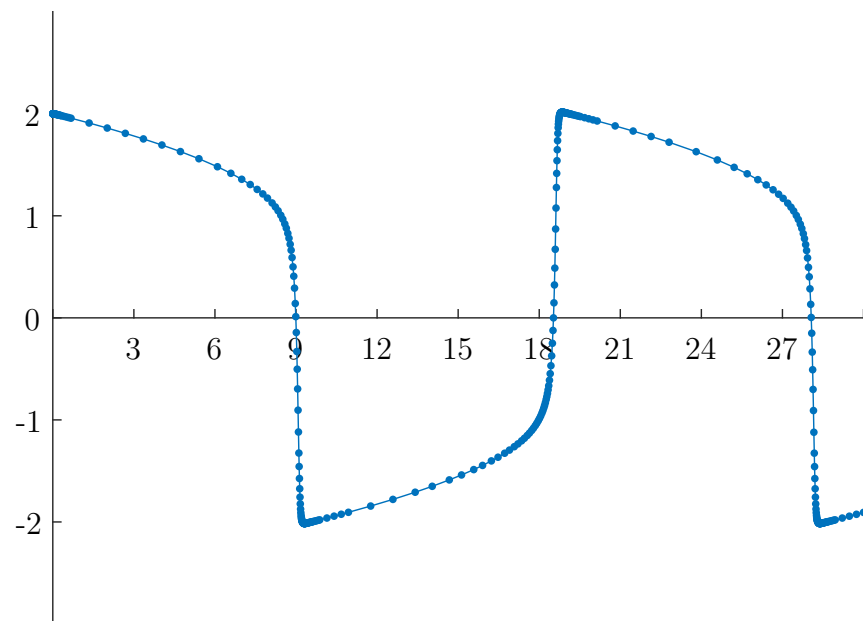


Figure 3.2: Van der Pol equation solution with `ode15s` ( $\mu = 10$ )

Which clearly uses less steps to pass through the stiff areas (a total of 326 with minimum stepsize 0.00014607).

If we where to solve it with a larger  $\mu$ , for example  $\mu = 1000$ , the number of steps needed using `ode45` is 5.495.393 which is too much compared to the 586 needed with `ode15s`.

**Definition.** The set of values of the stepsize  $h$  such that  $\lim_{h \rightarrow 0} y_n = 0$  is called the **absolute stability region** ( $\Omega$ ).

Let's see what happens with a general linear multistep method applied to a linear system:

Our method is

$$\underbrace{\sum_{j=0}^k \alpha_j y_{n+j}}_{\rho} = h \underbrace{\sum_{j=0}^k \beta_j f_{n+j}}_{\sigma}$$

We apply it to

$$y' = Ay$$

where we'll assume that the eigenvalues of  $A$  ( $\lambda_1, \dots, \lambda_n$ ) are all different (so that it's diagonalizable) and have negative real parts.

$y = e^{Ax}$  is the fundamental solution, and if  $y(0) = y_0$ , the solution is

$$y(x) = e^{Ax} y_0$$

with

$$\lim_{x \rightarrow \infty} e^{Ax} = \begin{pmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{pmatrix}$$

Now, with the system  $y' = Ay$ , our general method looks like

$$\sum_{j=0}^k \alpha_j y_{n+j} = h \sum_{j=0}^k \beta_j A y_{n+j}$$

Let's rewrite it:

$$\begin{aligned} & \sum_{j=0}^k (\alpha_j I - h \beta_j A) y_{n+j} = 0 \\ & \Downarrow \\ & \sum_{j=0}^k \left( \alpha_j I - h \beta_j \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} \right) y_{n+j} = 0 \\ & \Downarrow \\ & \sum_{j=0}^k (\alpha_j I - h \beta_j \lambda_i) y_{n+j} = 0 \quad (\forall i = 1, \dots, n) \end{aligned}$$

A finite difference equation looks like

$$a_k y_{n+k} + a_{k-1} y_{n+k-1} + \dots + a_0 y_n = 0$$

Which has roots  $r_1, \dots, r_k$ . So the solution is:

$$y_n = c_1 r_1^n + c_2 r_2^n + \dots + c_k r_k^n \xrightarrow{\text{if } \operatorname{Re}(r_k) < 0} 0$$

**Proposition 3.0.1.** Given

$$\hat{h} = h\lambda_i, \quad \rho(z) = \sum \alpha_j z^j, \quad \sigma(z) = \sum \beta_j z^j$$

the method is absolutely stable for values of  $\hat{h}$  such that the roots of  $\rho(z) - \hat{h}\sigma(z) = 0$  have negative real parts.

Let's see the one dimensional case with Euler's method:

**Example 3.0.2** (Stability region with Euler's method)

$$\left. \begin{array}{l} y_{n+1} = y_n + hf(y_n) \\ f(y) = \lambda y \end{array} \right\} \rightsquigarrow y_{n+1} = y_n + h\lambda y_n = (1 + h\lambda)y_n$$

$$\implies y_n = y_0 c(1 + h\lambda)^n \xrightarrow{\text{if } |1 + h\lambda| < 1} 0$$

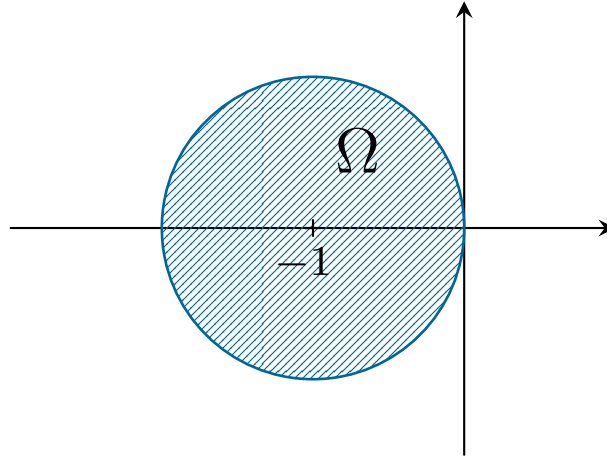


Figure 3.3: For all eigenvalues,  $\hat{h}$  must be in  $\Omega$  (stability region)

### Particular example

$$\begin{cases} y_1' = -2y_1 + y_2 + 2 \sin t \\ y_2' = 998y_1 - 999y_2 + 999(\cos t - \sin t) \end{cases}$$

This is a stiff system, with

$$J = \begin{pmatrix} -2 & 1 \\ 998 & -999 \end{pmatrix} \quad \text{Eigenvalues}(J) = \{-1000, -1\}$$

So the stepsize  $h$  needs to verify

$$|1 + h(-1000)| < 1$$

$$\implies h \leq \frac{1}{500}$$

Let's do it now with backwards Euler:

**Example 3.0.3** (Stability region with backwards Euler)

$$\begin{aligned} y_{n+1} &= y_n + \hat{h}y_{n+1} \rightsquigarrow (1 - \hat{h})y_{n+1} = y_n \rightsquigarrow y_{n+1} = \frac{1}{1 - \hat{h}}y_n \\ \implies \|1 - \hat{h}\| &= \|\hat{h} - 1\| > 1 \end{aligned}$$

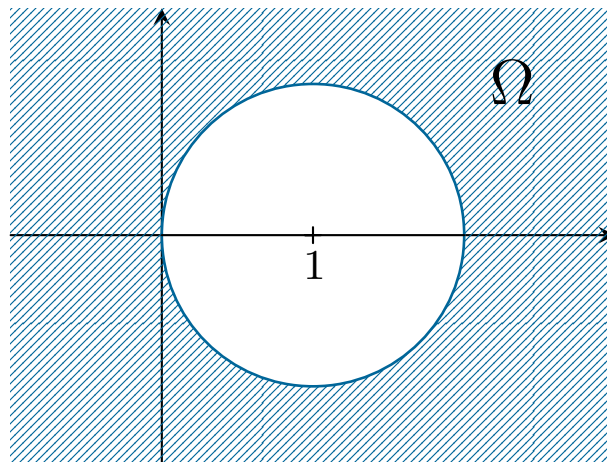


Figure 3.4: Stability region  $\Omega$  with backwards Euler

**Remark.** If all eigenvalues are  $< 1$ , we can take any  $h$ , as all  $\hat{h}$  are  $< 0$ .

**Example 3.0.4** (Dekker-Verwer p.7-8)

# PARTIAL DIFFERENTIAL EQUATIONS



## 4 Partial Differential Equations. Generalities on their solution

### 4.1 Finite differences

**Example 4.1.1** (1D Heat equation)

Our problem is:

$$\begin{cases} u_t - u_{xx} = f \\ u(x, 0) = u_0(x) \\ \left. \begin{array}{l} u(a, t) = u_a \\ u(b, t) = u_b \end{array} \right\} \end{cases} \quad \begin{array}{l} \leftarrow \text{Initial condition (IC)} \\ \leftarrow \text{Boundary conditions (BC)} \end{array}$$

With  $t \geq 0$ ,  $x \in [a, b]$

We discretize  $x$  and  $t$ :



Idea: We'll impose  $U_t(x_i, t^n) - U_{xx}(x_i, t^n) = f(x_i, t^n)$

### 4.1.1 Numerical derivatives



We'll use Taylor's expansion on  $f$  to find an approximation for the derivatives:

$$(1) \quad f_{i+1} = f_i + hf'_i + \frac{h^2}{2}f''_i + \frac{h^3}{3!}f'''_i + \dots$$

$$(2) \quad f_{i-1} = f_i - hf'_i + \frac{h^2}{2}f''_i - \frac{h^3}{3!}f'''_i + \dots$$

Approximation of the first derivative:

- Forward differences (1)

$$f'_i = \frac{f_{i+1} - f_i}{h} + \mathcal{O}(h)$$

- Backwards differences (2)

$$f'_i = \frac{f_i - f_{i-1}}{h} + \mathcal{O}(h)$$

- Centered differences (1)-(2)

$$f'_i = \frac{f_{i+1} - f_{i-1}}{2h} + \mathcal{O}(h^2)$$

**Remark.** With centered differences, it converges faster, but sometimes the other two methods are preferable for ease in computations.

### Approximation of the second derivative (1)+(2)

$$f_i'' = \frac{f_{i+1} - 2f_i + f_{i-1}}{h^2} + \mathcal{O}(h)$$

There are higher order approximations, but these are the standard, and are more than enough.

Now, we'll present some notation to differentiate between the problem we want to solve and the numerical problem.

#### **Definition.**

The **continuous problem** is 
$$\begin{cases} \mathcal{L}(u) = f & \text{in } \Omega, t > 0 \\ \mathcal{B}(u) = q & \text{in } \partial\Omega, t > 0 \\ u(x, 0) = u_0(x) \end{cases}$$

Where  $\mathcal{L}$  is any differential operator.

It's what we want to solve.

The **discrete problem** is 
$$\begin{cases} L(u) + \tau = f \\ B(u) + \bar{\tau} = q \\ u(x, 0) = u_0(x) \end{cases}$$

The **numerical problem** is 
$$\begin{cases} L(U) = f \\ B(U) = q \\ U(x, 0) = u_0(x) \end{cases}$$

Let's go back to our 1D heat equation:

The continuous problem is

$$\begin{cases} u_t - u_{xx} = f & t > 0, x \in (a, b) \\ u(a, t) = u_a \\ u(b, t) = u_b \\ u(x, 0) = u_0(x) \end{cases}$$

Notation:  $u(x_i, t^n) = u_i^n$

$x_i = a + i\Delta x \quad i = 0, \dots, N \quad x_0 = a, x_N = b \quad (\text{sometimes } \Delta x = h).$

$t^n = n\Delta t \quad n \geq 0$

We have the following explicit method:

#### 4.1.2 Forward Time Centered Space method (FTCS)

- Approximate  $u_t$  using forward differences (FT)
- Approximate  $u_{xx}$  using centered differences (CD)

Our discretized problem is

$$\frac{u_i^{n+1} - u_i^n}{\Delta t} + \mathcal{O}(\Delta t) - \frac{u_{i+1}^n - u_i^n + u_{i-1}^n}{\Delta x^2} + \mathcal{O}(\Delta x^2) = f_i^n$$

We define  $\tau = \mathcal{O}(\Delta t, \Delta x^2)$

Neglecting  $\tau$ , the numerical problem we want to solve is

$$\begin{cases} \frac{U_i^{n+1} - U_i^n}{\Delta t} - \frac{U_{i+1}^n - U_i^n + U_{i-1}^n}{\Delta x^2} = f_i^n \\ U_0^{n+1} = u_a(t^{n+1}) \\ U_N^{n+1} = u_b(t^{n+1}) \\ U_i^0 = u_0(x_i) \end{cases}$$

Now we define  $r$  as

$$r = \frac{\Delta t}{\Delta x^2}$$

So

$$U_i^{n+1} = rU_{i-1}^n + (1 - 2r)U_i^n + rU_{i+1}^n + \Delta t f_i^n \quad i = 1, \dots, N - 1$$

For  $i = 1$ ,

$$U_1^{n+1} = (1 - 2r)U_1^n + rU_2^n + \Delta t f_1^n + ru_a^n$$

And for  $i = N - 1$ ,

$$U_{N-1}^{n+1} = rU_{N-2}^n + (1 - 2r)U_{N-1}^n + \Delta t f_{N-1}^n + ru_b^n$$

Thus,

$$\overline{U^{n+1}} = \underset{\sim}{B} \overline{U^n} + \overline{F} + \overline{G}$$

Where

$$\begin{aligned} \overline{U^p} &= (U_1^p, U_2^p, \dots, U_{N-1}^p)^T & \overline{F} &= \Delta t (f_1^n, f_2^n, \dots, f_{N-1}^n)^T & \overline{G} &= (ru_a^n, 0, \dots, 0, ru_b^n)^T \\ \underset{\sim}{B} &= \begin{pmatrix} 1-2r & r & & & \\ r & 1-2r & r & & \\ & \ddots & \ddots & \ddots & \\ & & r & 1-2r & r \\ & & & r & 1-2r \end{pmatrix} \end{aligned}$$

...

## 5 Numerical Solution of PDEs with the Finite Elements Method

### 5.1 The Finite Elements Method (FEM)

#### 5.1.1 From the strong form to the weak form

Given the following problem in **strong form**:

$$\begin{cases} \nabla \cdot (\tilde{A} \nabla u) = -s & \text{in } \Omega \\ u = u_D & \text{on } \Gamma_D \\ \tilde{A} \frac{\partial u}{\partial n} = q & \text{on } \Gamma_N \end{cases}$$

Where

$$\begin{aligned} \Gamma_D \cup \Gamma_N &= \partial\Omega \\ \Gamma_D \cap \Gamma_N &= \emptyset \end{aligned}$$

$s$  is the source term (it's a function), but for the moment we'll assume that both  $s$  and  $\tilde{A}$  don't depend on  $u$ .

We'll derive the equivalent weak form of the problem using **weighted residuals**:

$$\int_{\Omega} \omega \nabla \cdot (\tilde{A} \nabla u) d\Omega = - \int_{\Omega} \omega s d\Omega \quad \forall \omega \text{ such that } \omega = 0 \text{ on } \Gamma_D$$

Where  $\omega$  is the **test function** (or **weighting function**).

We integrate by parts (in particular we use **Green's first identity**):

$$- \int_{\Omega} \nabla \omega \cdot (\tilde{A} \nabla u) d\Omega + \int_{\partial\Omega} \omega \tilde{A} \nabla u \cdot \bar{n} d\Gamma = - \int_{\Omega} \omega s d\Omega$$

Let's have a look at the term on  $\partial\Omega$

$$\int_{\partial\Omega} \omega \tilde{A} \nabla u \cdot \bar{n} d\Gamma = \int_{\Gamma_D} \overset{0 \text{ (by def)}}{\cancel{\omega \tilde{A} \nabla u \cdot \bar{n}}} d\Gamma + \int_{\Gamma_N} \omega \underbrace{\tilde{A} \nabla u \cdot \bar{n}}_q d\Gamma = \int_{\Gamma_N} \omega q d\Gamma$$

And with that, we get

#### Weak form

We have to find  $u$  such that  $u = u_D$  on  $\Gamma_D$  and

$$\int_{\Omega} \nabla \omega \cdot (\tilde{A} \nabla u) d\Omega = \int_{\Omega} \omega s d\Omega + \int_{\Gamma_N} \omega q d\Gamma \quad \forall \omega \text{ such that } \omega = 0 \text{ on } \Gamma_D$$

So

$$u \in \mathcal{H}_{\Gamma_D}^1(\Omega) = \{u \in \mathcal{H}^1(\Omega) \mid u = u_D \text{ on } \Gamma_D\}$$

and

$$\omega \in \mathcal{H}_{\Gamma_D,0}^1(\Omega) = \{u \in \mathcal{H}^1(\Omega) \mid u = 0 \text{ on } \Gamma_D\}$$

Where

$$\mathcal{H}^1(\Omega) = \left\{ f \in L^2(\Omega) \mid \frac{\partial f}{\partial x_i} \in L^2(\Omega) \quad \forall i \right\} \quad L^2(\Omega) = \left\{ f \mid \int_{\Omega} f^2 d\Omega < \infty \right\}$$

Now that we have the weak form, what we do is discretize the domain.

### 5.1.2 Discretisation

The mesh is defined by a set of points and triangles like:

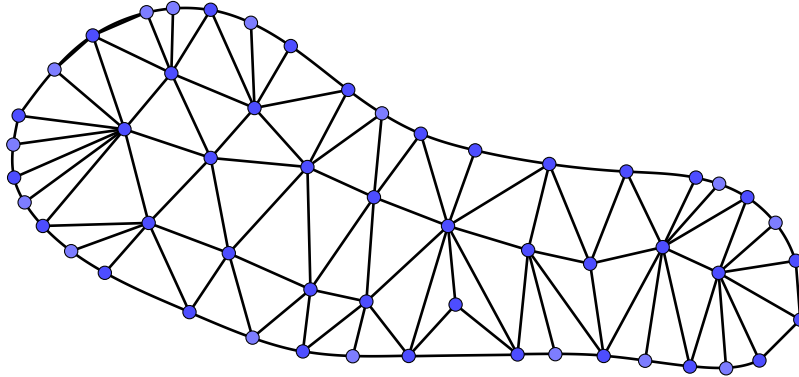


Figure 5.1: Mesh of a domain  $\Omega$

We want the nodal values of the solution noted  $u_i$  ( $i = 1, \dots, n_{\text{nodes}}$ ).

We'll approximate the solution  $u$  by:

$$u \simeq \sum_i N_i(x) u_i$$

**Definition.** We define the **shape function**  $N_i(x)$  as

$$N_i(x_j) = \delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

where  $x_i$  is a node on the mesh.

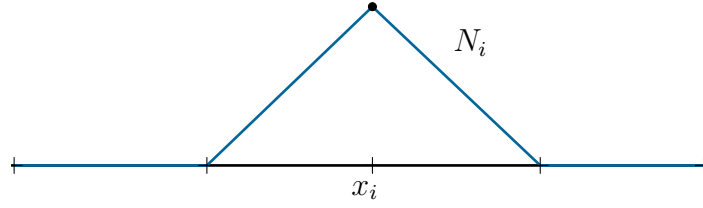


Figure 5.2:  $N_i$  are piecewise polynomials with compact support

Previously, we had to choose  $\omega$  as the weighting function. The common approach is with the *Galerkin approximation*:

$$\omega = N_i \quad i = 1, \dots, n_{\text{nodes}}$$

For every  $\omega$  chosen, we have an equation, so we choose as many  $\omega$ 's as unknowns.

Then

$$\begin{aligned} \int_{\Omega} \nabla N_i \cdot \tilde{A} \left( \sum_j \nabla N_j u_j \right) d\Omega &= \int_{\Omega} N_i s d\Omega + \int_{\Gamma_N} N_i q d\Gamma \\ \sum_j \underbrace{\int_{\Omega} \left( \nabla N_i \cdot \tilde{A} \nabla N_j \right) d\Omega}_{K_{ij}} u_j &= \underbrace{\int_{\Omega} N_i s d\Omega + \int_{\Gamma_N} N_i q d\Gamma}_{f_i} \end{aligned}$$

So we end up with the system

$$\sum_j K_{ij} u_j = f_i \quad i = 1, \dots, n_{\text{nodes}}$$

and in matrix form

$$\tilde{K} \bar{u} = \bar{f}$$

The **stiffness matrix**  $\tilde{K}$  is singular, so we need to impose the Dirichlet conditions:

$$\sum_{j=1}^{n_{\text{nodes}}} N_j(x) u_j = \sum_{j \in U} N_j(x) u_j + \sum_{j \in D} N_j(x) u_j \quad (u_j = u_D(x_j))$$

where  $\begin{cases} j \in D \text{ if } x_j \in \Gamma_D \\ j \in U \text{ if } x_j \notin \Gamma_D \end{cases}$



**Example 5.1.1**

We'll work with the following 1D problem:

$$\begin{cases} -u_{xx} = f \\ u(0) = 0, u(1) = 1 \end{cases}$$

with

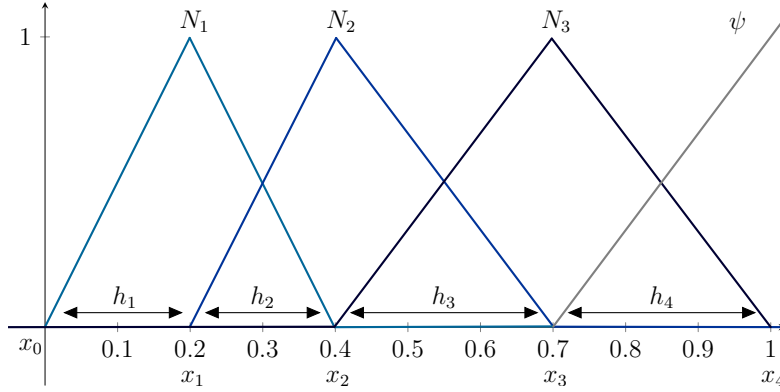


Figure 5.3

First we'll write the weak form to approximate  $u(x)$  by functions in  $\mathcal{H}_0^1(0, 1)$  to get

$$u(x) \simeq u^h = u_1 N_1(x) + u_2 N_2(x) + u_3 N_3(x) + \psi(x) \quad (5.1)$$

Let's proceed

$$-\int_0^1 u_{xx} \omega \, dx = \int_0^1 f \omega \, dx$$

We integrate the left hand side by parts:

$$-\int_0^1 u_{xx} \omega \, dx = -[\omega u_x]_0^1 + \int_0^1 u_x \omega_x \, dx$$

$$\omega \in \mathcal{H}_0^1(0, 1) \implies \omega(0) = \omega(1) = 0 \implies \omega u_x \Big|_0^1 = 0$$

**Weak form**

We have to find  $u$  such that  $u(0) = 0$ ,  $u(1) = 1$  and

$$\int_0^1 u_x \omega_x \, dx = \int_0^1 f \omega \, dx \quad \forall \omega \in \mathcal{H}_0^1(0, 1)$$

Let's rewrite  $\int_0^1 u_x \omega_x dx$  using our approximation of  $u$  (5.1):

$$\int_0^1 u_x \omega_x dx = \int_0^1 (u_1 N_1'(x) + u_2 N_2'(x) + u_3 N_3'(x) + \psi'(x)) \cdot N_i'(x) dx \quad i = 1, 2, 3$$

Thus

$$\sum_{j=1}^3 u_j \underbrace{\int_0^1 N_j'(x) \cdot N_i'(x) dx}_{K_{ij}} = \underbrace{\int_0^1 f N_i(x) dx}_{f_i} - \underbrace{\int_0^1 \psi'(x) N_i'(x) dx}_{f_i} \quad i = 1, 2, 3$$

Let's find the explicit form of the  $N_i'$ s with the discretisation in Figure 5.3:

$$N_1(x) = \begin{cases} 5x & 0 \leq x \leq 0.2 \\ -5(x - 0.4) & 0.2 \leq x \leq 0.4 \\ 0 & \text{otherwise} \end{cases}$$

$$N_2(x) = \begin{cases} 5(x - 0.2) & 0.2 \leq x \leq 0.4 \\ -\frac{10}{3}(x - 0.7) & 0.4 \leq x \leq 0.7 \\ 0 & \text{otherwise} \end{cases}$$

$$N_3(x) = \begin{cases} \frac{10}{3}(x - 0.4) & 0.4 \leq x \leq 0.7 \\ -\frac{10}{3}(x - 1) & 0.7 \leq x \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

Now we can find the  $K_{ij}$ 's:

$$K_{11} = \int_0^1 N_1'(x) \cdot N_1'(x) dx = \int_0^{0.2} 5 \cdot 5 dx + \int_{0.2}^{0.4} (-5) \cdot (-5) dx = 10$$

$$K_{12} = K_{21} = \int_{0.2}^{0.4} (-5) \cdot 5 dx = -5$$

$$K_{13} = K_{31} = 0$$

$$K_{22} = \int_{0.2}^{0.4} 5 \cdot 5 dx + \int_{0.4}^{0.7} \frac{-10}{3} \cdot \frac{-10}{3} dx = \frac{25}{3}$$

$$K_{23} = K_{32} = \int_{0.4}^{0.7} \frac{-10}{3} \cdot \frac{10}{3} dx = -\frac{10}{3}$$

$$K_{33} = \int_{0.4}^{0.7} \frac{10}{3} \cdot \frac{10}{3} dx + \int_{0.7}^1 \frac{-10}{3} \cdot \frac{-10}{3} dx = \frac{20}{3}$$

So

$$\underset{\sim}{K} = \begin{pmatrix} 10 & -5 & 0 \\ -5 & \frac{25}{3} & -\frac{10}{3} \\ 0 & -\frac{10}{3} & \frac{20}{3} \end{pmatrix}$$

and the system we have to solve is:

$$\begin{pmatrix} 10 & -5 & 0 \\ -5 & \frac{25}{3} & -\frac{10}{3} \\ 0 & -\frac{10}{3} & \frac{20}{3} \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ f_3 \end{pmatrix}$$

Let's find a general way to compute  $\underset{\sim}{K}$  with discretisations like the one in Figure 5.3:

$$K_{11} = \int_0^1 N'_1 N'_1 dx = \int_{x_0}^{x_1} \left( \frac{1}{h_1} \right)^2 dx + \int_{x_1}^{x_2} \left( -\frac{1}{h_2} \right)^2 dx = \frac{1}{h_1} + \frac{1}{h_2}$$

$$K_{12} = K_{21} = \int_0^1 N'_1 N'_2 dx = \int_{x_1}^{x_2} -\frac{1}{h_2} \cdot \frac{1}{h_2} dx = -\frac{1}{h_2}$$

$\vdots$

$$\Rightarrow \underset{\sim}{K} = \begin{pmatrix} \frac{1}{h_1} + \frac{1}{h_2} & -\frac{1}{h_2} & 0 \\ -\frac{1}{h_2} & \frac{1}{h_2} + \frac{1}{h_3} & -\frac{1}{h_3} \\ 0 & -\frac{1}{h_3} & \frac{1}{h_3} + \frac{1}{h_4} \end{pmatrix}$$

The solution needs to verify the Dirichlet conditions, so we'll impose

$$\psi = 0 \cdot N_0 + 1 \cdot N_4 = N_4$$

Assuming  $f$  constant, the system we have to solve is

$$\begin{pmatrix} \frac{1}{h_1} + \frac{1}{h_2} & -\frac{1}{h_2} & 0 \\ -\frac{1}{h_2} & \frac{1}{h_2} + \frac{1}{h_3} & -\frac{1}{h_3} \\ 0 & -\frac{1}{h_3} & \frac{1}{h_3} + \frac{1}{h_4} \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} = \begin{pmatrix} f \cdot \frac{h_1+h_2}{2} \\ f \cdot \frac{h_2+h_3}{2} \\ f \cdot \frac{h_3+h_4}{2} \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ -\left(\frac{-1}{h_4}\right) \end{pmatrix} \quad (5.2)$$

$\Downarrow$

$$\begin{pmatrix} 10 & -5 & 0 \\ -5 & \frac{25}{3} & -\frac{10}{3} \\ 0 & -\frac{10}{3} & \frac{20}{3} \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} = f \cdot \begin{pmatrix} 0.2 \\ 0.25 \\ 0.3 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ \frac{10}{3} \end{pmatrix}$$

Taking for example  $f = -2$ ,

$$\bar{u} = (0.04, 0.16, 0.49)$$

And the approximation of the solution is

$$u^h(x) = \sum_{i=0}^4 u_i N_i(x) = 0.04 N_1(x) + 0.16 N_2(x) + 0.49 N_3(x) + N_4(x)$$

Of course, this is a toy problem. If things get more complicated, we analyze the problem locally:

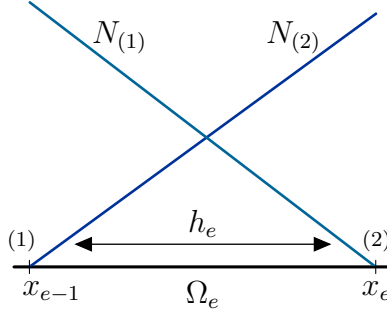


Figure 5.4

With

$$\int_0^1 N'_i N'_j dx = \sum_{e=1}^4 \int_{\Omega_e} N'_i N'_j dx = \sum_{e=1}^4 \int_{x_{e-1}}^{x_e} N'_i N'_j dx$$

$$N'_{(1)} = -\frac{1}{h_e} \quad N'_{(2)} = \frac{1}{h_e}$$

$$\int_{\Omega_e} (N'_{(1)})^2 dx = \frac{1}{h_e} \quad \int_{\Omega_e} (N'_{(2)})^2 dx = \frac{1}{h_e} \quad \int_{\Omega_e} N'_{(1)} N'_{(2)} dx = -\frac{1}{h_e}$$

$$\int_{\Omega_e} N_{(1)} f dx = f \frac{h_e}{2} \quad \int_{\Omega_e} N_{(2)} f dx = f \frac{h_e}{2}$$

And so, we define the **element matrix** as

$$\begin{pmatrix} \frac{1}{h_e} & -\frac{1}{h_e} \\ -\frac{1}{h_e} & \frac{1}{h_e} \end{pmatrix}$$

Then, we assemble our system as follows:

$$\begin{pmatrix} \frac{1}{h_1} & \frac{-1}{h_1} & & & \\ \frac{-1}{h_1} & \frac{1}{h_1} + \frac{1}{h_2} & \frac{-1}{h_2} & & \\ & \frac{-1}{h_2} & \frac{1}{h_2} + \frac{1}{h_3} & \frac{-1}{h_3} & \\ & & \frac{-1}{h_3} & \frac{1}{h_3} + \frac{1}{h_4} & \frac{-1}{h_4} \\ & & & \frac{-1}{h_4} & \frac{1}{h_4} \end{pmatrix} \begin{pmatrix} u_0 \\ u_1 \\ u_2 \\ u_3 \\ u_4 \end{pmatrix} = f \cdot \begin{pmatrix} \frac{h_1}{2} \\ \frac{h_1}{2} + \frac{h_2}{2} \\ \frac{h_2}{2} + \frac{h_3}{2} \\ \frac{h_3}{2} + \frac{h_4}{2} \\ \frac{h_4}{2} \end{pmatrix}$$

To get the previous system (5.2), we cross out the first and last rows, and move the highlighted terms to the right hand side:

$$\begin{pmatrix} \frac{1}{h_1} & \frac{-1}{h_1} & & & \\ \cancel{\frac{-1}{h_1}} & \frac{1}{h_1} + \frac{1}{h_2} & \frac{-1}{h_2} & & \\ & \frac{-1}{h_2} & \frac{1}{h_2} + \frac{1}{h_3} & \frac{-1}{h_3} & \\ & & \frac{-1}{h_3} & \frac{1}{h_3} + \frac{1}{h_4} & \cancel{\frac{-1}{h_4}} \\ & & & \frac{-1}{h_4} & \frac{1}{h_4} \end{pmatrix} \begin{pmatrix} u_0 \\ u_1 \\ u_2 \\ u_3 \\ u_4 \end{pmatrix} = f \cdot \begin{pmatrix} \frac{h_1}{2} \\ \frac{h_1}{2} + \frac{h_2}{2} \\ \frac{h_2}{2} + \frac{h_3}{2} \\ \frac{h_3}{2} + \frac{h_4}{2} \\ \frac{h_4}{2} \end{pmatrix} + \begin{pmatrix} 0 \\ u_0 \cdot \frac{1}{h_1} \\ 0 \\ u_4 \cdot \frac{1}{h_4} \\ 0 \end{pmatrix}$$

$\Downarrow$

$$\begin{pmatrix} \frac{1}{h_1} + \frac{1}{h_2} & \frac{-1}{h_2} & 0 \\ \frac{-1}{h_2} & \frac{1}{h_2} + \frac{1}{h_3} & \frac{-1}{h_3} \\ 0 & \frac{-1}{h_3} & \frac{1}{h_3} + \frac{1}{h_4} \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} = f \cdot \begin{pmatrix} \frac{h_1}{2} + \frac{h_2}{2} \\ \frac{h_2}{2} + \frac{h_3}{2} \\ \frac{h_3}{2} + \frac{h_4}{2} \end{pmatrix} + \begin{pmatrix} u_0 \cdot \frac{1}{h_1} \\ 0 \\ u_4 \cdot \frac{1}{h_4} \end{pmatrix}$$

And imposing the boundary conditions ( $u_0 = 0$ ,  $u_4 = 1$ ) we obtain the exact same system (5.2).

Now, we'll see with an example how to assemble the stiffness matrix in 2D:

### Example 5.1.2

## 5.2 FEM analysis

**Remark.** In this section we'll only see the 1D case, but the analysis works for higher dimensions.

Let's work with the following problem:

$$\begin{cases} \Delta u = -s & \text{in } \Omega \\ u = 0 & \text{on } \partial\Omega \end{cases} \quad (5.3)$$

### Weak form

We have to find  $u \in \mathcal{H}_0^1(\Omega)$  such that

$$\underbrace{\int_{\Omega} \nabla \omega \nabla u \, d\Omega}_{a(\omega, u)} = \underbrace{\int_{\Omega} s \omega \, d\Omega}_{l(\omega)} \quad \forall \omega \in \mathcal{H}_0^1(\Omega)$$

### Discrete problem

Find  $u^h \in V^h$ , ( $V^h \subseteq \mathcal{H}_0^1(\Omega)$ ) such that

$$a(\omega^h, u^h) = l(\omega^h) \quad \forall \omega^h \in V^h$$

We want to compute the error  $e = u - u^h$

$$\|u - u^h\| \leq C \cdot \inf_{v^h \in V^h} \|u - v^h\| \leq C \cdot \|u - \Pi u\|$$

where  $\Pi u$  is the projection of  $u$ . In FEM it's the *interpolating polynomial*.

**Lemma 5.2.1** (Lax-Milgram).

Let  $V$  be a Hilbert space. Let

$$a : V \times V \rightarrow \mathbb{R}$$

a continuous, coercive and bilinear form.

Given the weak form:

$$\text{Find } u \in V \text{ such that } a(\omega, u) = l(\omega) \quad \forall \omega \in V \quad (5.4)$$

Then (5.4) has a unique solution and

$$\|u\|_V \leq \frac{1}{\alpha} \|l\|_V$$

Here,

$a$  is **continuous** if  $\exists M > 0$  such that  $|a(\omega, u)| \leq M \|u\|_V \|\omega\|_V$

$a$  is **coercive** if  $\exists \alpha > 0$  such that  $a(u, u) \geq \alpha \|u\|_V^2$

If moreover  $a(\cdot, \cdot)$  is symmetric, then (5.4) is equivalent to

$$\text{Find } u \in V \text{ such that } J(u) = \min_{v \in V} J(v), \text{ with } J(v) = \frac{1}{2} a(v, v) + l(v)$$

**Example 5.2.1**

Let's see that for problem (5.3), the hypothesis of Lax-Milgram's lemma holds:

Of course,  $a$  is a bilinear form. Then, we also want to see that  $a(\omega, u) = \int_{\Omega} \nabla \omega \nabla u \, d\Omega$  is continuous and coercive.

We work in  $V \subseteq \mathcal{H}_0^1(\Omega)$  with the norm  $\|u\|_V := \|\nabla u\|$ , and  $\langle u, v \rangle = \int_{\Omega} uv$ . Then

$$|a(\omega, u)| = \left| \int_{\Omega} \nabla \omega \nabla u \, d\Omega \right| = |\langle \nabla \omega, \nabla u \rangle| \underset{\text{Cauchy-Schwarz}}{\leq} \|\nabla \omega\| \cdot \|\nabla u\| = \|\omega\|_V \cdot \|u\|_V \implies a \text{ continuous}$$

$$a(u, u) = \int_{\Omega} (\nabla u)^2 \, d\Omega = \|u\|_V^2. \text{ In particular } a(u, u) \geq \|u\|_V^2 \implies a \text{ coercive}$$

**Lemma 5.2.2** (Céa).

The Galerkin method is **strongly consistent**. That is

$$a(v^h, \underbrace{u - u^h}_e) = 0 \quad \forall v^h \in V^h$$

**Remark.** We are assuming that  $a$  verifies the hypothesis of Lax-Milgram's lemma (5.2.1)

*Proof.* Evaluating the weak form (5.4) on a subspace  $V^h \subset V$ , it holds

$$a(v^h, u) = l(v^h) \tag{5.5}$$

$$a(v^h, u^h) = l(v^h) \tag{5.6}$$

Subtracting (5.5) - (5.6) we obtain the desired result

$$a(v^h, u - u^h) = 0$$

Where we used that  $a$  is bilinear. □

**Corollary 5.2.1.** In the same conditions as before, we have the following bound:

$$\|u - u^h\|_V \leq \frac{M}{\alpha} \|u - v^h\|_V \quad \forall v^h \in V^h$$

*Proof.*

$$\begin{aligned} a(e, e) &= a(u - u^h, u - u^h) \stackrel{(*)}{=} \\ &= a(u - u^h, u - v^h) + a(u - u^h, \cancel{v^h - u^h}) \stackrel{0 \text{ (Céa's Lemma)}}{\leq} \stackrel{(**)}{\leq} \\ &\leq M \|u - u^h\|_V \|u - v^h\|_V \end{aligned}$$

In (\*) we added and subtracted  $v^h$  and used that  $a$  is a bilinear form, and in (\*\*) we used that  $a$  is continuous.

Now, as  $a$  is coercive, we also have

$$\begin{aligned} a(e, e) &= a(u - u^h, u - u^h) \leq \alpha \|u - u^h\|_V^2 \\ \implies \alpha \|u - u^h\|_V^2 &\leq M \|u - u^h\|_V \|u - v^h\|_V \\ \implies \|u - u^h\|_V &\leq \frac{M}{\alpha} \|u - v^h\|_V \end{aligned}$$

□

This result means that the Galerkin approximation  $u^h$  is the best we can obtain (up to a constant depending only on the form  $a$ ).

Now our objective is to find an error bound for the FEM solution. We'll use a bound for the interpolation error (We'll do it with 1D FEM for simplicity).

We have the domain  $\Omega = (a, b)$ . Given  $x_0 = a, x_1, x_2, \dots, x_{N-1}, x_N = b$  such that

$$x_0 = a < x_1 < x_2 < \dots < x_{N-1} < x_N = b$$

we define  $\tau_h$  as the partition of  $(a, b)$  in elements  $K_j = (x_{j-1}, x_j)$  (with  $j = 1, 2, \dots, N$ )

and

$$X_h^p = \left\{ v^h \in \mathcal{C}^0(\Omega) : v^h|_{K_j} \in \mathbb{P}^p(K_j) \right\} \quad \varphi_i(x_j) = \delta_{ij}$$

where  $\mathbb{P}^p$  are polynomials of degree  $p$ .

**Theorem 5.2.1.** Let  $v \in \mathcal{H}^{p+1}(I)$  and  $\Pi_h^p v \in X_h^p$  its interpolating function. Then

$$|v - \Pi_h^p v|_{\mathcal{H}^k(I)} \leq C_{p,k} \cdot h^{p+1-k} \cdot |v|_{\mathcal{H}^{p+1}(I)} \quad k = 0, 1, \dots$$

In our case, for  $p = 1$ ,

$$\begin{aligned} \Pi_h^1 v &\in X_h^1(\Omega) \\ \Pi_h^1 v(x_j) &= v(x_j) \quad j = 0, \dots, N \\ \varphi_j(x) &= \begin{cases} \frac{x - x_{j-1}}{x_j - x_{j-1}} & \text{if } x \in K_j = (x_{j-1}, x_j) \\ \frac{x - x_j}{x_{j+1} - x_j} & \text{if } x \in K_{j+1} = (x_j, x_{j+1}) \\ 0 & \text{otherwise} \end{cases} \quad (\text{Lagrange basis functions}) \end{aligned}$$

$$\Pi_h^1 v = \sum_{j=0}^N v(x_j) \varphi_j(x)$$



*Proof.* ( $p = 1$ )

$$e = v - \Pi_h^1 v$$

Note that  $e(x_j) = 0 \xRightarrow[\text{Rolle}]{} e'(\xi_j) = 0$  with  $\xi_j \in K_j$

In  $K_j$ ,  $e'(x) = \int_{\xi_j}^x e''(s) ds = \int_{\xi_j}^x v''(s) ds$ , and so

$$|e'(x)| \leq \int_{x_{j-1}}^{x_j} |v''(s)| ds \leq \underbrace{\left( \int_{x_{j-1}}^{x_j} 1 ds \right)^{\frac{1}{2}}}_{h^{\frac{1}{2}}} \left( \int_{x_{j-1}}^{x_j} |v''(s)|^2 ds \right)^{\frac{1}{2}}$$

We square both sides

$$|e'(x)|^2 \leq h \int_{x_{j-1}}^{x_j} |v''(s)|^2 ds$$

and integrate in  $K_j$

$$\int_{x_{j-1}}^{x_j} |e'(x)|^2 dx \leq h \cdot \underset{\substack{\text{Comes from integrating} \\ \int_{x_{j-1}}^{x_j} \int_{x_{j-1}}^{x_j} ( ) = \int_{x_{j-1}}^{x_j} ( ) \cdot \underbrace{(x_j - x_{j-1})}_h}}{h} \int_{x_{j-1}}^{x_j} |v''(s)|^2 ds \quad (*)$$

Then we have  $e(x) = \int_{x_{j-1}}^x e'(s) ds$  for  $x \in K_j$ . Thus

$$|e(x)| \leq \int_{x_{j-1}}^{x_j} |e'(s)| ds \leq h^{\frac{1}{2}} \left( \int_{x_{j-1}}^{x_j} |e'(s)|^2 ds \right)^{\frac{1}{2}} \stackrel{(*)}{\leq} h^{\frac{1}{2}} \cdot h \cdot \left( \int_{x_{j-1}}^{x_j} |v''(s)|^2 ds \right)^{\frac{1}{2}}$$

and again, if we square and integrate both sides, we end up with

$$\int_{x_{j-1}}^{x_j} |e(x)|^2 dx \leq h^4 \int_{x_{j-1}}^{x_j} |v''(s)|^2 ds$$

To summarize, we have

$$\begin{aligned} \int_{K_j} |e(x)|^2 dx &\leq h^4 \int_{K_j} |v''(s)|^2 ds \\ \int_{K_j} |e'(x)|^2 dx &\leq h^2 \int_{K_j} |v''(s)|^2 ds \end{aligned}$$

and if we sum for all the elements,

$$\|e(x)\|_{\mathcal{L}_2}^2 \leq h^4 |v|_{\mathcal{H}^2}^2 \implies \|e\|_{\mathcal{L}_2} \leq C_1 h^2$$

$$|e(x)|_{\mathcal{H}^1}^2 \leq h^2 |v|_{\mathcal{H}^2}^2 \implies |e|_{\mathcal{H}^1} \leq C_2 h$$

□

**Remark.** In  $\mathcal{H}^p(\Omega)$  we have

- Norm:  $\|v\|_{\mathcal{H}^p} = \sum_{k=0}^p |v|_{\mathcal{H}^k}$
- Seminorm:  $|v|_{\mathcal{H}^k}^2 = \int_{\Omega} \left( \frac{D^k v}{Dx^k} \right)^2 d\Omega$

**Corollary 5.2.2.**

$$\begin{aligned} \|u - u^h\|_{\mathcal{H}^p(\Omega)} &\leq \frac{M}{\alpha} \inf_{v^h \in \mathcal{H}_0^p(\Omega)} \|u - v^h\|_{\mathcal{H}_0^p(\Omega)} \leq \frac{M}{\alpha} Ch^p \\ &\quad \uparrow \\ &\quad u^h = \Pi^h u \\ \|u - u^h\|_{\mathcal{L}_{\epsilon}} &\leq Kh^{p+1} \end{aligned}$$