MATH 536 - Numerical Solution of Partial Differential Equations

February 13, 2020

ii

# Contents

1	Systems of Linear Equations						
	1.1	Direct Methods – Gaussian Elimination					
		1.1.1 Forward Elimination					
	1.2	Gaussian Elimination with Pivoting					
	1.3	Special Cases					
		1.3.1 Cholesky Decomposition					
		1.3.2 Thomas Algorithm					
	1.4	Matrix Norms					
	1.5	Error Estimates					
	1.6	Iterative Method Preliminaries					
	1.7	Iterative Methods					
		1.7.1 Matrix Splitting Methods					
		1.7.2 Optimization-based Methods 12					
2	Solı	utions to Partial Differential Equations 23					
-	2.1	Classification of Partial Differential Equations 23					
	2.1	2.1.1 First Order linear PDFs					
		2.1.2 Second Order PDE					
	2.2	Difference Operators					
		2.2.1 Poisson Equation					
		2.2.2 Neumann Boundary Conditions					
	2.3	Consistency and Convergence					
	2.4	Advection Equation					
	2.5	Von Neumann Stability Analysis					
	2.6	Sufficient Conditions for Convergence					
	2.7	Parabolic PDEs – Heat Equation					
		2.7.1 FC Scheme					
		2.7.2 BC scheme					
		2.7.3 Crank-Nicolson Scheme					
		2.7.4 Leapfrog Scheme					
		2.7.5 DuFort-Frankel Scheme					
	2.8	Advection-Diffusion Equation					
		2.8.1 FC Scheme					
		2.8.2 Upwinding Scheme (FB)					
9	Tota	reduction to Finite Floments (12)					
ა	2 1	Weighted Residual Methods 43					
	0.1	$\begin{array}{cccccccccccccccccccccccccccccccccccc$					
		3.1.1 Suburt spaces					
		2.1.2 Weighted Residual Formulations					
		5.1.5 Conocation Methods					

## CONTENTS

3.2	Weak Methods
3.3	Finite Element Method (FEM) 49
3.4	Gaussian Quadrature
3.5	Error Estimates
3.6	Optimal Error Estimates
	3.6.1 Other Boundary Conditions
3.7	Transient Problems
3.8	Finite Elements in Two Dimensions
	3.8.1 Finite Element Basis Functions
	3.8.2 Discrete Galerkin Formulation

# Chapter 1

# Systems of Linear Equations

# 1.1 Direct Methods – Gaussian Elimination

### 1.1.1 Forward Elimination

(k-th step):

$$\begin{aligned} a_{ij}^{[k]} &= a_{ij}^{[k-1]} - a_{k-1,j}^{[k-1]} \ell_{i,k-1}^{[k-1]}, \qquad i = k, \dots, n, \quad j = k-1, \dots, n \\ \ell_{i,k-1}^{[k-1]} &= \frac{a_{i,k-1}^{[k-1]}}{a_{k-1,k-1}^{[k-1]}}, \qquad k = 2, \dots, n \quad i = k, \dots, n \\ b_i^{[k]} &= b_i^{[k-1]} - b_{k-1}^{[k-1]} \ell_{i,k-1}^{[k-1]}, \qquad i = k, \dots, n \end{aligned}$$

In n-1 steps:

$$\begin{bmatrix} a_{11}^{[1]} & a_{12}^{[1]} & \dots & a_{1n}^{[1]} \\ 0 & a_{22}^{[2]} & \dots & a_{2n}^{[2]} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & a_{nn}^{[n]} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1^{[1]} \\ b_2^{[2]} \\ \vdots \\ b_n^{[n]} \end{bmatrix}$$

or

Ux = b

#### **Operations Count**

$$\sum_{j=1}^{n-1} \left[ \underbrace{(n-j)}_{\text{divisions}} + 2\underbrace{(n-j)(n-j+1)}_{\text{multiply and sum}} \right] = \frac{2n^3}{6} + \frac{n^2}{2} - \frac{7n}{6}$$

Where we have used

$$\sum_{j=1}^{m} j = \frac{m(m+1)}{2}, \quad \sum_{j=1}^{m} j^2 = \frac{m(m+1)(2m+1)}{6}$$

#### **Backward Substitution**

$$x_{i} = \frac{1}{a_{ii}^{[i]}} \left[ b_{i}^{[i]} - \sum_{j=i+1}^{n} a_{ij}^{[i]} x_{j} \right]$$
1

**Operation Count** 

$$\sum_{i=1}^{n} \left[ \underbrace{\frac{2(n-i)}{\text{Multiplication and Subtraction}} + \underbrace{1}_{\text{division}} \right] = n^2$$

#### Matrix Representations

Define:

$$L_{k} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & \ddots & 0 & \dots & 0 \\ 0 & \dots & 1 & 0 & \dots & 0 \\ 0 & \dots & -\ell_{k+1,k}^{[k]} & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & -\ell_{nk}^{[k]} & 0 & \dots & 1 \end{bmatrix}$$

Note that:

$$L_k^{-1} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & \ddots & 0 & \dots & 0 \\ 0 & \dots & 1 & 0 & \dots & 0 \\ 0 & \dots & \ell_{k+1,k}^{[k]} & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \ell_{nk}^{[k]} & 0 & \dots & 1 \end{bmatrix}$$

Lemma 1.1.1

$$A^{[2]} = L_1 A^{[1]}, \quad b^{[2]} = L_1 b^{[1]} \quad \left(A^{[1]} = A, b^{[1]} = b\right)$$

Therefore,

or

$$A = L_1^{-1} L_2^{-1} \cdots L_{n-1}^{-1} U = L U$$

 $A^{[n]} = U = L_{n-1}L_{n-2}\cdots L_1A$ 

#### Proposition 1.1.2

$$L = L_1^{-1} \cdots L_{n-1}^{-1}$$

is a lower triangular matrix with unit diagonal.

**Definition** A *Permutation matrix* is a matrix that has exactly one nonzero entry in each row and column, and this entry is equal to 1. The elementary permutation matrix (EPM) is defined as a matrix produced from the identity matrix by interchanging exactly one pair k, m of its columns. We denote it as  $P_{km} = P_k$ .

For example:

$$P_{23} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

**Proposition 1.1.3** If  $A \in \mathbb{R}^{n \times n}$  and  $P_{k,m} \in \mathbb{R}^{n \times n}$  is EPM, then  $P_{k,m}A$  differs from A by an interchange of rows k and m.  $AP_{k,m}$  differs from A by an interchange of columns k and m.

# 1.2 Gaussian Elimination with Pivoting

$$A^{[n]} = U = L_{n-1}P_{n-1}\dots L_1P_1A$$

Where U is upper triangular,  $L_k$  is defined as before, and  $P_k = P_{km}$  with  $m \ge k$ . Equivalently,

$$A = \underbrace{P_1 L_1^{-1} P_2 L_2^{-1} \cdots P_{n-1} L_{n-1}^{-1}}_{L^*} U$$

**Lemma 1.2.1** If  $i \ge j > k$  and  $P_j = P_{ji}$  then  $P_j L_k^{-1} P_j$  is produced form  $L_k^{-1}$  by interchanging the *j*-th and *i*-th entry in the *k*-th column.

#### Proof

$$\begin{bmatrix} \ddots & & & & \\ 1 & & & \\ \vdots & \ddots & & \\ \ell_{jk}^{[k]} & 1 & & \\ \vdots & \ddots & & \\ \ell_{ik}^{[k]} & & 1 & \\ \vdots & & \ddots & \\ \ell_{ik}^{[k]} & & 1 & \\ \vdots & & & \ddots \end{bmatrix} \xrightarrow{\underline{P_j L_k^{-1}}} \begin{bmatrix} \ddots & & & & \\ 1 & & & \\ \vdots & \ddots & & \\ \ell_{ik}^{[k]} & 1 & 0 & \\ \vdots & & & \ddots \end{bmatrix} \xrightarrow{\underline{P_j L_k^{-1} P_j}} \begin{bmatrix} \ddots & & & & \\ 1 & & & \\ \vdots & \ddots & & \\ \ell_{ik}^{[k]} & 1 & 0 & \\ \vdots & & & \ddots \end{bmatrix}$$

**Theorem 1.2.2** Consider  $P = P_{n-1} \dots P_1$ . Then P is a permutation matrix and  $PA = PL^*U = LU$ , with L being a lower triangular matrix with unit diagonal.

Proof

$$PL^{*} = PP_{1}L_{1}^{-1}P_{2}L_{2}^{-1}\cdots P_{n-1}L_{n-1}^{-1} = P_{n-1}P_{n-2}\cdots \underbrace{P_{1}P_{1}}_{I}L_{1}^{-1}P_{2}L_{2}^{-1}\cdots P_{n-1}L_{n-1}^{-1} = P_{n-1}P_{n-2}\cdots P_{2}L_{1}^{-1}P_{2}\underbrace{P_{3}\cdots P_{n-1}P_{n-1}\cdots P_{3}}_{I}L_{2}^{-1}\cdots P_{n-1}L_{n-1}^{-1} = L$$

$$L_{1*}^{-1}P_{n-1}\cdots P_{3}L_{2}^{-1}P_{3}\cdots P_{n-1}L_{n-1}^{-1} = \cdots = L_{1*}^{-1}L_{2*}^{-1}\cdots L_{n-2*}^{-1}L_{n-1}^{-1} = L$$

**Definition** If  $A \in \mathbb{R}^{n \times n}$  and  $1 \le k \le n$  then

$$\hat{A}_k = \begin{bmatrix} a_{11} \dots a_{1k} \\ a_{21} \dots a_{2k} \\ \vdots & \ddots & \vdots \\ a_{k1} \dots & a_{kk} \end{bmatrix}$$

is called the  $k^{\text{th}}$  principle sub-matrix of A

**Theorem 1.2.3** The pivot entries  $a_{kk}^{[k]}$ , k = 1, 2, ..., n-1 are nonzero if and only if  $\hat{A}_k$  are non-singular for k = 1, ..., n-1 (note that  $\hat{A}_n = A$  and we assume that A is nonsingular since otherwise the linear system is ill-posed).

**Proof** (i) Suppose first that all  $a_{kk}^{[k]} \neq 0, k = 1, ..., n-1$ . Then since  $\hat{A}_k = \hat{L}_k \hat{U}_k$  it follows that  $det(\hat{A}_k) = det(\hat{U}_k) = a_{11}^{[k]} \dots a_{kk}^{[k]} \neq 0$ 

(ii) Now suppose that  $det(\hat{A}_k) \neq 0, k = 1, \dots, n-1$ . Then, using a induction argument, we have:

(a)  $a_{11}^{[1]} = \hat{A}_1$  and therefore  $a_{11}^{[1]} \neq 0$ . (b) Assume that  $a_{11}^{[1]}, \ldots, a_{kk}^{[k]} \neq 0$ . If  $\hat{A}_{k+1}^{[k+1]}$  is the k+1 principle sub-matrix of  $A^{[k+1]}$  then it is easy to see that  $\hat{A}_{k+1}^{[k+1]} = (\hat{L}_k)_{k+1} \ldots (\hat{L}_1)_{k+1} \hat{A}_{k+1}^{[1]}$  where we again assume that  $(\hat{L}_m)_{k+1}, m = 1, \ldots, k$  are the k+1 principle sub-matrices of  $L_m$  and  $\hat{A}_{k+1}^{[1]}$  is the k+1 principle sub-matrix of  $A^{[1]} = A$ . Therefore,  $det(\hat{A}_{k+1}^{[k+1]}) = det(\hat{A}_{k+1}) \neq 0$ . But  $\hat{A}_{k+1}^{[k+1]}$  is upper triangular i.e.

$$\hat{A}_{k+1}^{[k+1]} = \begin{bmatrix} a_{11}^{[1]} a_{12}^{[1]} \dots a_{1,k+1}^{[1]} \\ 0 & a_{22}^{[2]} \dots & a_{2,k+1}^{[2]} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & a_{k+1,k+1}^{[k+1]} \end{bmatrix}$$

and then  $det(\hat{A}_{k+1}^{[k+1]}) = a_{11}^{[1]} \dots a_{k+1,k+1}^{[k+1]}$ . This, together with the induction hypothesis yields that  $a_{k+1,k+1}^{[k+1]} \neq 0$  which completes the proof.

**Definition** A matrix is *strictly diagonally dominant* if

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|$$

**Corollary 1.2.4** If a matrix is strictly diagonally dominant then no pivoting is necessary.

Follows from theorem 1.2.3 and the following theorem due to Gershgorin:

**Theorem 1.2.5** The spectrum of  $A \in \mathbb{C}^{n \times n}$ , S(A), is enclosed in the set:

$$\left(\bigcup_{i=1}^n D_i\right),\,$$

where:

$$D_i = \left\{ z \in \mathbb{C} : |z - a_{ii}| \le \sum_{j \ne i} |a_{ij}| \right\}$$

**Definition** Matrix A is symmetric positive definite (spd) if

- 1.  $A = A^T$
- 2.  $v^T A v \ge 0$  for all  $v \in \mathbb{R}^n$
- 3.  $v^T A v = 0$  if and only if  $v \equiv 0$

Corollary 1.2.6 If a matrix is spd then no pivoting is necessary.

Follows from the following theorem.

**Theorem 1.2.7** If  $A \in \mathbb{R}^{n \times n}$  is symmetric the following statements are equivalent

- 1. A is positive definite,
- 2. S(A) contains only positive real numbers, and
- 3. Every principle sub-matrix is positive definite.

# 1.3 Special Cases

#### 1.3.1 Cholesky Decomposition

**Theorem 1.3.1** If  $A \in \mathbb{R}^{n \times n}$  is spd then there exists a lower triangular matrix C such that  $CC^T = A$ . Furthermore, the diagonal entries of C are positive.

**Proof** Use induction on *n*:

$$n = 1$$
 :  $C = [c_{11}] = \sqrt{a_{11}}$ 

Assume the theorem holds for all matrices in  $\mathbb{R}^{n \times n}$ .

$$A = \begin{bmatrix} A_n & a \\ a^T & a_{n+1,n+1} \end{bmatrix}$$

where  $a_{n+1,n+1} > 0$ . By the induction hypothesis:

$$A_n = C_n C_n^T$$

Now try to find [X] such that

$$CC^{T} = \begin{bmatrix} C_{n} & 0\\ X^{T} & c \end{bmatrix} \begin{bmatrix} C_{n}^{T} & X\\ 0 & c \end{bmatrix} = \begin{bmatrix} A_{n} & a\\ a^{T} & a_{n+1,n+1} \end{bmatrix}$$

or such that

$$C_n X = a \quad \Rightarrow \quad X = C_n^{-1} a$$

and

$$X^T X + c^2 = a_{n+1,n+1} \quad \Rightarrow \quad c = \sqrt{a_{n+1,n_1} - X^T X}$$

But, is c > 0?

$$X^{T}X = \left[C_{n}^{-1}a\right]^{T}C_{n}^{-1}a = a^{T}\left[C_{n}^{-1}\right]^{T}C_{n}^{-1}a = a^{T}A_{n}^{-1}a$$

therefore

$$a_{n+1,n+1} - a^T A_n^{-1} a = a_{n+1,n+1} - X^T X$$

Then define  $x := \left[ \left[ A_n^{-1} a \right]^T, -1 \right]^T \neq 0$  so that

$$x^{T}Ax = a_{n+1,n+1} - a^{T}A_{n}^{-1}a > 0$$

#### Cholesky Algorithm

1. 
$$c_{11} \leftarrow \sqrt{a_{11}}$$
  
2. For  $i = 2, 3, ..., n$   
3.  $c_{i1} \leftarrow a_{i1}/c_{11}$   
4. End  $i$   
5. For  $j = 2, 3, ..., n - 1$   
6.

$$c_{jj} \leftarrow \sqrt{a_{jj} - \sum_{k=1}^{j-1} c_{jk}^2}$$

7. For i = j + 1, ..., n

8.

$$c_{ij} \leftarrow \frac{1}{c_{jj}} \left( a_{ij} - \sum_{k=1}^{j-1} c_{ik} c_{jk} \right)$$

9. End i

10. End j

11.

$$c_{nn} \leftarrow \sqrt{a_{nn} - \sum_{k=1}^{n-1} c_{nk}^2}$$

12. End

#### 1.3.2 Thomas Algorithm

**Definition** A *band matrix* is of the form

$a_{11}$	$a_{1,2}$		$a_{1,l}$	0					0 -
$a_{21}$	·	·		·	·				÷
:	·	·	·		۰.	·			÷
$a_{k,1}$		$a_{k,k-1}$	$a_{kk}$	$a_{k,k+1}$		$a_{k,k+l-1}$	0		0
0	·		·	·	·		·.	·	÷
:	·	·		·	·	·		·	0
:		·	·		۰.	·	·		$a_{n-l+1,n}$
:			·	·		·	·	·	÷
:				·	·		·	·	$a_{n-1,n}$
0	• • • • •		• • • • •		0	$a_{n,n-k+1}$		$a_{n,n-1}$	$a_{n,n}$

A band matrix may be stored compactly as an  $l + k - 1 \times n$  matrix of the form:



#### **Thomas Algorithm**

For the linear system

$$\begin{bmatrix} b_1 & c_1 & 0 & \dots & 0 \\ a_2 & b_2 & c_2 & 0 & \dots & 0 \\ 0 & a_3 & b_3 & c_3 & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \ddots & 0 \\ \vdots & & & \ddots & a_{n-1} & b_{n-1} & c_{n-1} \\ 0 & \dots & 0 & a_n & b_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_n \end{bmatrix}$$

6

1.

2.

$$b_{j}^{[f]} = b_{j} - a_{j} \frac{c_{j-1}^{[f]}}{b_{j-1}^{[f]}} \quad \text{and} \quad d_{j}^{[f]} = d_{j} - a_{j} \frac{d_{j-1}^{[f]}}{b_{j-1}^{[f]}}, \quad \text{for} \quad j = 2, \dots, n$$
$$x_{j-1} = d_{j-1}^{[f]} - b_{j}^{[f]} x_{j}, \quad \text{for} \quad j = n, \dots, 2$$

# 1.4 Matrix Norms

#### **Definition** If $A \in \mathbb{R}^{n \times n}$ then the mapping $\mathbb{R}^{n \times n} \to \mathbb{R}$ is called a norm of A := ||A|| if and only if:

- 1.  $||A|| \ge 0$  and ||A|| = 0 if and only if A = 0,
- 2.  $||\lambda A|| = |\lambda| ||A||$  for all  $\lambda \in \mathbb{R}$ ,
- 3.  $||A + B|| \le ||A|| + ||B||.$

If in addition,  $||AB|| \leq ||A|| ||B||$  the norm is called multiplicative. In the sequel, we make use of multiplicative norms only, and therefore, the notion of a norm presumes a multiplicative norm.

**Definition** If  $A \in \mathbb{R}^{n \times n}$ , then for a vector norm  $|| \cdot || : \mathbb{R}^n \to \mathbb{R}$ ,

$$||A|| := \sup_{||x|| \neq 0} \frac{||Ax||}{||x||}$$

is called a *subordinate* matrix norm.

**Theorem 1.4.1** Let  $A \in \mathbb{R}^{n \times n}$ . Then,

$$||A||_{\infty} := \sup_{||x||_{\infty} \neq 0} \frac{||Ax||_{\infty}}{||x||_{\infty}} = \max_{1 \le i \le n} \sum_{j=1}^{n} |a_{ij}|$$

2.

1.

$$||A||_{1} := \sup_{||x||_{1} \neq 0} \frac{||Ax||_{1}}{||x||_{1}} = \max_{1 \le j \le n} \sum_{i=1}^{n} |a_{ij}|$$

3.

$$||A||_{2} := \sup_{||x||_{2} \neq 0} \frac{||Ax||_{2}}{||x||_{2}} = \sqrt{\rho(A^{T}A)}$$

Proof

$$||Ax||_{\infty} = \max_{1 \le i \le n} \left| \sum_{j=1}^{n} a_{ij} x_j \right| \le \left( \max_{1 \le i \le n} \sum_{j=1}^{n} |a_{ij}| \right) \left( \max_{1 \le j \le n} |x_j| \right) = M ||x||_{\infty}$$

 $||A||_{\infty} \le M$ 

therefore,

Let I be the index for which

$$\max_{1 \le i \le n} \sum_{j=1}^{n} |a_{ij}| = M$$

is achieved. Choose  $x = \{x_j\}_{j=1}^n$ 

$$x_j := \begin{cases} a_{Ij} / |a_{Ij}| & : & a_{Ij} \neq 0 \\ 0 & : & a_{ij} = 0. \end{cases}$$

Then  $||x||_{\infty} = 1$  and  $||Ax||_{\infty} = M$ . Therefore,

$$||A||_{\infty} = M$$

# **1.5** Error Estimates

Consider two vectors  $x_1$  and  $x_2$  where

$$Ax_1 = b$$
 and  $Ax_2 = b + r$ 

then

$$x_1 = A^{-1}b$$
 and  $x_2 = A^{-1}(b+r)$ 

 $||x_2 - x_1|| = ||A^{-1}r||$ 

and it follows immediately that

Therefore,

$$\frac{||x_2 - x_1||}{||A|| \, ||x_1||} \le \frac{||x_2 - x_1||}{||Ax_1||} = \frac{||A^{-1}r||}{||b||} \le \frac{||A^{-1}|| \, ||r||}{||b||}$$

$$\frac{||x_2 - x_1||}{||x_1||} \le \underbrace{||A|| \, ||A^{-1}||}_{c(A)} \frac{||r||}{||b||}$$

This provides means of estimating the error introduced in the solution of a linear system due to roundoff errors or uncertainty in the right-hand-side or the matrix of the system. The number c(A) is called the condition number of A with respect to the norm  $\|.\|$ .

# **1.6** Iterative Method Preliminaries

**Definition** If  $x_k \in \mathbb{R}^n$  then the sequence  $\{x_k\}_{k=1}^{\infty}$  converges to  $x \in \mathbb{R}^n$  in a norm  $|| \cdot ||$  i.e.

$$\lim_{k \to \infty} x_k = x$$

if

$$\lim_{k \to \infty} ||x - x_k|| = 0.$$

**Definition** If  $A_k \in \mathbb{R}^{n \times n}$  then the sequence  $\{A_k\}_{k=1}^{\infty}$  converges to  $A \in \mathbb{R}^{n \times n}$  in a norm  $|| \cdot ||$  i.e.

$$\lim_{k \to \infty} A_k = A$$

if

$$\lim_{k \to \infty} ||A - A_k|| = 0.$$

Consider the matrix power series

$$\sum_{k=0}^{\infty} a_k A^k. \tag{1.1}$$

It is convergent if and only if

$$\lim_{K \to \infty} \sum_{k=0}^{K} a_k A^k = f(A) \,,$$

where f(A) is some matrix with a finite norm. In the following we will also need to consider the numerical power series

$$f(\lambda) \sim \sum_{k=0}^{\infty} a_k \lambda^k.$$
(1.2)

#### 1.7. ITERATIVE METHODS

**Theorem 1.6.1** The matrix power series (1.1) is convergent if for all  $\lambda_i \in S(A) |\lambda_i| < \rho$ , where  $\rho$  is the radius of convergence of (1.2) and S(A) is the spectrum of A. It is divergent if  $|\lambda_i| > \rho$  for some i.

**Proof** From Jordan's theorem, there exists nonsingular C such that A can be transformed as follows:

$$B = C^{-1}AC, \ B = \operatorname{diag}(B_i), \ B_i = \begin{vmatrix} \lambda_i & 1 & 0 & \dots & 0 \\ 0 & \lambda_i & 1 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & 1 \\ 0 & \dots & 0 & \lambda_i \end{vmatrix}$$

where  $B_i \in \mathbb{R}^{n_i \times n_i}$  and  $n_i$  is the multiplicity of the eigenvalue  $\lambda_i$ . Then,  $\sum_{k=0}^{\infty} a_k A^k$  is convergent if and only if

$$\sum_{k=0}^{\infty} a_k B^k = \sum_{k=0}^{\infty} a_k C^{-1} A^k C = C^{-1} \left( \sum_{k=0}^{\infty} a_k A^k \right) C$$

is convergent as well. Note that  $B^k = \text{diag}(B_i^k)$  since B is a block-diagonal matrix. Now, consider the powers  $B_i^k$  of the *i*-th block of B. It can be shown by induction that:

$$B_{i}^{2} = \begin{bmatrix} \lambda_{i}^{2} & 2\lambda_{i} & 1 & \dots & 0\\ 0 & \lambda_{i}^{2} & 2\lambda_{i} & \dots & 0\\ \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & \lambda_{i}^{2} \end{bmatrix}, \dots, \quad B_{i}^{k} = \begin{bmatrix} \lambda_{i}^{k} & \binom{k}{1}\lambda_{i}^{k-1} & \dots & \binom{k}{n_{i}-1}\lambda_{i}^{k-n_{i}+1}\\ 0 & \lambda_{i}^{k} & \binom{k}{1}\lambda_{i}^{k-1} & \dots & \binom{k}{n_{i}-2}\lambda_{i}^{k-n_{i}+2}\\ \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & \dots & \lambda_{i}^{k} \end{bmatrix}$$

so that the *i*-th block of the *m*-th partial sum of  $\sum_{k=0}^{\infty} a_k B^k$  is

$$f_m(B_i) = \sum_{k=0}^m a_k B_i^k = \begin{bmatrix} f_m(\lambda_i) & \frac{1}{1!} f'_m(\lambda_i) & \dots & \frac{1}{(n_i-1)!} f_m^{(n_i-1)}(\lambda_i) \\ 0 & f_m(\lambda_i) & \dots & \frac{1}{(n_i-2)!} f_m^{(n_i-2)}(\lambda_i) \\ \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & f_m(\lambda_i) . \end{bmatrix}$$

Therefore,  $f_m(B)$  is convergent if and only if

$$f_m\left(\lambda_i\right) = \sum_{k=0}^m a_k \lambda_i^k$$

(and therefore  $f_m^{(l)}(\lambda_i)$ ) is convergent for i = 1, ..., n and  $l = 1, ..., n_i$ .  $f_m(\lambda_i)$  is convergent if  $|\lambda_i|$  is smaller than  $\rho$ , the radius of convergence of (1.2).

#### Corollary 1.6.2

$$f(A) = I + A + A^2 + \dots + A^m + \dots$$

is convergent if and only if  $|\lambda_k| < 1$  for all k = 1, ..., n and  $A \in \mathbb{R}^{n \times n}$ .

**Corollary 1.6.3** Since  $|\lambda_k| \leq ||A||$  for any subordinate norm  $||\cdot||$  then f(A) is convergent if ||A|| < 1.

#### **1.7** Iterative Methods

#### 1.7.1 Matrix Splitting Methods

Consider

Ax = b

If A = B + C, so that  $det(C) \neq 0$ , then

$$Ax = b \quad \Leftrightarrow \quad x = x + C^{-1} (b - Ax)$$

and in this form we may solve the problem iteratively:

$$x^{k+1} = x^k + C^{-1} \left( b - A x^k \right) \tag{1.3}$$

or equivalently:

$$Cx^{k+1} = b - Bx^k. (1.4)$$

The last form clarifies why such methods are called matrix-splitting methods. Note, that C can be any nonsingular matrix but not every choice is a good choice, of course. The essence of various iterative methods is in the choice of the matrix C that allows for the fast computation of a good approximation to the solution. If C is the identity matrix, the corresponding iteration is called simple (or Richardson) iteration. It is seldom a good choice. The form (1.3) suggests that the closer C is to A, the faster the convergence. On the other hand, the iteration requires the solution of a linear system with C and therefore, C should be such that this solution requires much less resources than the solution of the original system. These are very contradictory requirements and the choice of C depends on the properties of A and some other contraints like computer resources.

**Theorem 1.7.1** The iteration (1.3) is convergent if and only if  $\rho(I - C^{-1}A) < 1$ .

Proof

$$x^{k+1} = \underbrace{(I - C^{-1}A)}_{D} x^{k} + C^{-1}b = D^{2}x^{k-1} + (I + D)C^{-1}b = \dots$$
$$= D^{k+1}x^{0} + (I + D + D^{2} + \dots + D^{k})C^{-1}b$$

 $I + D + \ldots + D^k + \ldots$  is convergent iff  $\rho(I - C^{-1}A) < 1$ . Also, if it is convergent then

$$D^k \xrightarrow[k \to \infty]{} 0$$

**Definition**  $R(D) = -\log_{10} \rho(D)$  is called the *convergence rate* of the iteration.

Jacobi's Method

$$C = \text{diag}(A)$$
$$D = I - C^{-1}A = \begin{bmatrix} 0 & -\frac{a_{12}}{a_{11}} & -\frac{a_{13}}{a_{11}} & \cdots & -\frac{a_{1n}}{a_{11}} \\ \vdots & \ddots & \vdots \\ -\frac{a_{n1}}{a_{nn}} & -\frac{a_{n2}}{a_{nn}} & -\frac{a_{n3}}{a_{nn}} & \cdots & 0 \end{bmatrix}$$

**Theorem 1.7.2** If A is strictly diagonally dominant then Jacobi is convergent.

Proof

$$||D||_{\infty} = \max_{i} \sum_{j \neq i} \left| \frac{a_{ij}}{a_{ii}} \right| < 1.$$

Then for any eigenvalue  $\lambda_i$  of D we have:

$$|\lambda_i| \le ||D||_{\infty} < 1 \qquad .$$

10

#### 1.7. ITERATIVE METHODS

#### Gauss-Seidel (GS) Iteration

$$C = L = \begin{bmatrix} a_{11} & 0\\ \vdots & \ddots\\ a_{1n} & \dots & a_{nn} \end{bmatrix}$$
$$B = U = A - L$$
$$x^{k+1} = x^k + C^{-1} (b - Ax^k)$$

**Theorem 1.7.3** If A is strictly diagonally dominant the Gauss-Seidel iteration is convergent.

**Proof** The formula for the *i*-th component of the k + 1 iterate  $x^{k+1}$  is:

$$x_i^{k+1} = -\sum_{j < i} \frac{a_{ij}}{a_{ii}} x_j^{k+1} - \sum_{j > i} \frac{a_{ij}}{a_{ii}} x_j^k + \frac{b_i}{a_{ii}}$$

The exact solution clearly is a fixed point of the iteration i.e.

$$x_i = -\sum_{j < i} \frac{a_{ij}}{a_{ii}} x_j - \sum_{j > i} \frac{a_{ij}}{a_{ii}} x_j + \frac{b_i}{a_{ii}}$$

Therefore the error of the k + 1 iteration,  $\epsilon^{k+1} = x^{k+1} - x$ , must satisfy:

$$\epsilon_i^{k+1} = -\sum_{j < i} \frac{a_{ij}}{a_{ii}} \epsilon_j^{k+1} - \sum_{j > i} \frac{a_{ij}}{a_{ii}} \epsilon_j^k$$

then

$$\left|\epsilon_{i}^{k+1}\right| \leq \sum_{j < i} \left|\frac{a_{ij}}{a_{ii}}\right| \left|\epsilon_{j}^{k+1}\right| + \sum_{j > i} \left|\frac{a_{ij}}{a_{ii}}\right| \left|\epsilon_{j}^{k}\right| \tag{1.5}$$

Using induction on i, it is straightforward to show that:

$$\left|\epsilon_{i}^{k+1}\right| \leq \left(\sum_{j \neq i} \left|\frac{a_{ij}}{a_{ii}}\right|\right) \left|\left|\epsilon^{k}\right|\right|_{\infty} < \left|\left|\epsilon^{k}\right|\right|_{\infty}.$$

Indeed, assuming that  $|\epsilon_j^{k+1}| < ||\epsilon^k||_{\infty}$ ,  $\forall j < i$ , and substituting this into (1.5) yields that  $|\epsilon_i^{k+1}| < ||\epsilon^k||_{\infty}$ . Therefore,

$$\left|\left|\epsilon^{k+1}\right|\right|_{\infty} < \left|\left|\epsilon^{k}\right|\right|_{\infty}$$

so that the error must strictly decrease with each iteration. Then the error must go to zero because the exact solution is the unique fixed point of the iteration (why is it unique?).

#### Successive Over-relaxation (SOR) Method

Suppose  $\omega \in \mathbb{R}$  and consider the iteration:

$$x^{k+1} = (\omega^{-1}D + L)^{-1} (b - ((1 - \omega^{-1}) D + U) x^{k})$$

Here  $C = L + \frac{1}{\omega}D$ , D = diag(A), and being L the lower triangular part of A, not including the main diagonal. Then clearly A = C + B, where  $B = (1 - \frac{1}{\omega})D + U$ , U being the strictly upper triangular part of A (with zeros on the main diagonal). This iteration is similar to the GS iteration but has a free parameter that allows to better control the convergence rate.

**Theorem 1.7.4** (Ostrowski-Reich) If  $A \in \mathbb{R}^{n \times n}$  is spd then SOR is convergent iff  $0 < \omega < 2$ .

#### 1.7.2 Optimization-based Methods

From now on, until the end of this section we assume that the matrix A is s.p.d.

**Theorem 1.7.5** If  $A \in \mathbb{R}^{n \times n}$  is spd then u is a solution to the linear system Au = b if and only if u minimizes the function:

$$F\left(v\right) = \frac{1}{2}v^{T}Av - b^{T}v$$

**Proof** Let u be the solution of the linear system i.e. Au = b. Then for any  $v \in \mathbb{R}$  we have that:

$$F(v) - F(u) = \frac{1}{2}v^{T}Av - b^{T}v - \frac{1}{2}u^{T}Au + b^{T}u$$
  
=  $\frac{1}{2}v^{T}Av - u^{T}Av + \frac{1}{2}u^{T}Au = \frac{1}{2}(v - u)^{T}A(v - u) \ge 0$ 

then

$$\frac{1}{2}(v-u)^{T}A(v-u) = 0 \quad \Leftrightarrow \quad v = u$$

so that  $F(v) \ge F(u)$  and

$$F(v) = F(u) \quad \Leftrightarrow \quad v = u$$

**Definition** If A is spd then for any two vectors  $u, v \in \mathbb{R}^n$ ,  $u^T A v$  defines the *energy inner product* induced by A denoted by:

$$\langle u, v \rangle_A = u^T A v.$$

Note that in the rest of the notes at some occasions  $(u, v)_A$  is also used to denote the energy inner product.

#### Steepest Descent Method

Note that any iteration for finding a solution to a linear system Au = b can be written in the form:

$$u^{k+1} = u^k + \alpha_k p_k$$

where  $\alpha_k \in \mathbb{R}$  is called iteration step,  $p_k \in \mathbb{R}^n$  is called the search direction. For example, in the case of matrix splitting methods the search direction is chosen to be  $p_k = C^{-1}r_k$ , with  $r_k = b - Au^k$  being called the residual of the iterate  $u^k$ , and  $\alpha_k = 1$ . The step does not need to be constant and the gradient-type methods choose it at each iteration step so that this choice minimizes the function F(v). The basic gradient-type iterative algorithm therefore can be written as:

1. Choose an initial search direction

$$p_0 = r_0 = b - Au_0$$

- 2. For  $k = 1, 2, 3, \dots$  do:
  - (a) Find  $\alpha_k \in \mathbb{R}$  such that

 $u^{k+1} = u^k + \alpha_k p_k$ 

minimizes F over the line  $u^k + \alpha p_k$ 

(b) New iterate

 $u^{k+1} = u^k + \alpha_k p_k$ 

(c) The new residual is

 $r_{k+1} = b - Au^{k+1}$ 

- (d) Find the new search direction  $p_{k+1}$ .
- (e) Let k = k + 1 and repeat item 2 until convergence.

#### 1.7. ITERATIVE METHODS

To find  $\alpha_k$  we minimize F along the search direction  $p_k$ :

$$\frac{d}{d\alpha_k}F\left(u^k + \alpha_k p_k\right) = \nabla F\left(u^k + \alpha_k p_k\right) \cdot \frac{d}{d\alpha_k}\left(u^k + \alpha_k p_k\right)$$
$$= \left[A\left(u^k + \alpha_k p_k\right) - b\right]^T p_k = \left[Au^k - b\right]^T p_k + \alpha_k p_k^T A p_k$$
$$= -r_k^T p_k + \alpha_k p_k^T A p_k = 0$$

so that

$$\alpha_k = \frac{r_k^T p_k}{p_k^T A p_k}$$

Note that for this choice of  $\alpha_k$  we have:

$$\left[\underbrace{Au^{k+1}-b}_{-r_{k+1}}\right]^T p_k = \left[A\underbrace{\left(u^k + \alpha_k p_k\right)}_{u^{k+1}} - b\right]^T p_k = -r_k^T p_k + \alpha_k p_k^T A p_k = 0$$
(1.6)

$$-r_{k+1}^{T}p_{k} = \nabla F(u_{k+1})p_{k} = 0$$
(1.8)

In the steepest descent method we choose

$$p_k = r_k = b - Au^k.$$

or

Because this is the direction of the negative gradient of F(v), this function is clearly non-increasing in this direction. Unfortunately, this choice of search direction is not optimal in some sense that will be elucidated in the next section.

#### Conjugate Gradient Method (CGM)

**Definition**  $u^k$  is optimal with respect to direction  $p \neq 0$  iff  $F(u^k) \leq F(u^k + \lambda p)$  for all  $\lambda \in \mathbb{R}$ .

**Lemma 1.7.6**  $u^k$  is optimal with respect to p iff

$$p^T r_k = 0$$

**Proof**  $u^k$  is optimal with respect to p iff  $F(u^k + \lambda p)$  has a minimum at  $\lambda = 0$ , i.e.

$$\frac{\partial F}{\partial \lambda} \left( u^k + \lambda p \right) \Big|_{\lambda=0} = p^T \underbrace{\left( A u^k - b \right)}_{-r_k} + \lambda p^T A p \Big|_{\lambda=0} = 0.$$
(1.9)

(Note that  $\frac{\partial^2 F}{\partial \lambda^2} = p^T A p > 0.$ ) (1.9) is obviously satisfied iff  $p^T r_k = 0$ .

We know that the choice of  $\alpha_k$  guarantees that  $u_{k+1}$  is optimal with respect to  $p_k$  and this applies in particular to the steepest descent method for which  $p_k = r_k$ . However,  $u^{k+2}$  may not be optimal with respect to  $p_k$  i.e. the next iteration can undo some of the minimization work of the previous iteration. Indeed:

$$u^{k+2} = u^{k+1} + \alpha_{k+1}r_{k+1}$$

which gives

$$r_{k+2} = r_{k+1} - \alpha_{k+1}Ap_{k+1} \Rightarrow p_k^T r_{k+2} = \underbrace{p_k^T r_{k+1}}_{=0} - \alpha_{k+1}p_k^T Ap_{k+1} \quad \Rightarrow \quad p_k^T r_{k+2} = -\alpha_{k+1}p_k^T Ap_{k+1}.$$

For the steepest descent method  $p_m = r_m$  and therefore, since  $r_k^T r_{k+1} = 0$  (see (1.8)),  $r_k^T A r_{k+1} \neq 0$ , unless A = cI.

To find a different search direction  $p_{k+1}$  that maintains the optimality of  $u^{k+2}$  w.r.t.  $p_k$ , we need that:

$$0 = p_k^T r_{k+2} = \alpha_{k+1} p_k^T A p_{k+1} = 0$$

For the conjugate gradient method we require that  $u^{k+2}$  is also optimal with respect to  $p_k$  i.e. we choose the search direction  $p_{k+1}$  from this condition.

**Definition** If  $u, v \in \mathbb{R}^n$  are such that  $\langle u, v \rangle_A = 0$  for some s.p.d. A then u, v are called A-conjugate.

This means that  $p_{k+1}$  should be *A*-conjugate to  $p_k$ . We search for  $p_{k+1}$  in the form  $p_{k+1} = r_{k+1} + \beta_k p_k$ , for some  $\beta_k \in \mathbb{R}$ , and then from the condition  $p_k^T A p_{k+1} = 0$  we easily obtain that  $\beta_k = -\frac{\langle r_{k+1}, p_k \rangle_A}{||p_k||_A^2}$ 

**Algorithm** (Basic CGM procedure) If  $A \in \mathbb{R}^{n \times n}$  is spd and  $u^0$  is an initial guess:

1. 
$$r_{0} \leftarrow b - Au^{0}$$
  
2.  $p_{0} = r_{0}$   
3. For  $k = 1, 2, ..., m$ :  
4.  $\alpha_{k-1} \leftarrow r_{k-1}^{T} p_{k-1} / ||p_{k-1}||_{A}^{2}$   
5.  $u^{k} \leftarrow u^{k-1} + \alpha_{k-1} p_{k-1}$   
6.  $r_{k} \leftarrow r_{k-1} - \alpha_{k-1} Ap_{k-1}$   
7.  $p_{k} \leftarrow r_{k} - \frac{\langle r_{k}, p_{k-1} \rangle_{A}}{||p_{k-1}||_{A}^{2}} p_{k-1}$ 

8. Next k

9. End

As it will be demonstrated below, if all arithmetic operations are exact, the algorithm computes the exact solution in a finite number of steps (finite termination property).

1

Later, we will show that the error decreases as:

$$\left|\left|u-u^{k}\right|\right|_{A}=\left|\left|\epsilon^{k}\right|\right|_{A}\leq2\left|\left|\epsilon_{0}\right|\right|_{A}\left[\frac{\sqrt{\operatorname{c}\left(A\right)}-1}{\sqrt{\operatorname{c}\left(A\right)}+1}\right]^{k}$$

If  $c(A) \gg 1$  then the convergence is slow. Therefore, the algorithm is often modified by formally multiplying the original system by a matrix that has a spectrum close to the spectrum of A and performing the CG method on the modified system which has the same solution as the original one. This process is called preconditioning.

#### Preconditioning

Instead of Au = b we solve  $\tilde{A}\tilde{u} = \tilde{b}$  where

$$\tilde{A} = (B^{-1})^T A B^{-1}, \quad \tilde{u} = B u, \quad \tilde{b} = (B^{-1})^T b$$

Since the preconditioner  $B^T B$  must be an approximation to the matrix A, a natural choice for B should be an approximation to the square root of A. Since A is s.p.d. its square root is also s.p.d. and therefore we assume that  $B^T = B$ . The preconditioner  $B^2$  must satisfy somewhat contradicting requirements since on one hand it

#### 1.7. ITERATIVE METHODS

should be a good approximation of A and on the other, it should be much easier to solve a linear system with a matrix  $B^2$  than with A. The second requirement follows from the fact that, as it will become clear below, on each iteration of the CG method with preconditioning, we will need to solve one system with  $B^2$ .

For the system Au = b we know that

$$u: F(u) = \min_{v \in \mathbb{R}^n} F(v)$$

Now, if  $\tilde{v} = Bv$  then

$$\begin{split} F\left(v\right) &= \frac{1}{2}v^{T}Av - b^{T}v = \frac{1}{2}\left[B^{-1}\tilde{v}\right]^{T}A\left[B^{-1}\tilde{v}\right] - b^{T}\left(B^{-1}\tilde{v}\right) \\ &= \frac{1}{2}\tilde{v}^{T}\underbrace{B^{-1}AB^{-1}}_{\tilde{A}}\tilde{v} - \underbrace{\left(B^{-1}b\right)^{T}}_{\tilde{b}^{T}}\tilde{v} = \tilde{F}\left(\tilde{v}\right) \end{split}$$

The preconditioned method should compute the new iterate from:  $\tilde{u}^k = \tilde{u}^{k-1} + \tilde{\alpha}_{k-1}\tilde{p}_{k-1}$ , with  $\tilde{\alpha}_{k-1} = \frac{\tilde{r}_{k-1}^T\tilde{p}_{k-1}}{||\tilde{p}_{k-1}||_{\tilde{A}}}$ . The search direction for the preconditioned method should be computed as  $\tilde{p}_k = \tilde{\beta}_{k-1}\tilde{p}_{k-1} + \tilde{r}_k$  where the residual  $\tilde{r}_k$  is given by  $\tilde{r}_k = \tilde{r}_{k-1} - \tilde{\alpha}_{k-1}\tilde{A}\tilde{p}_{k-1}$ , and  $\tilde{\beta}_k = -\frac{\langle \tilde{r}_{k+1}, \tilde{p}_k \rangle_{\tilde{A}}}{||\tilde{p}_k||_{\tilde{A}}^2}$ . These formulae are not convenient for practical computations since they require to somehow compute the inverse of the preconditioner,  $(B^2)^{-1}$ . Fortunately, it appears that the preconditioned method can be implemented almost exactly as the original CGM modifying only the computation of the search direction to:

$$p_{k} = \underbrace{\left(B^{2}\right)^{-1} r_{k}}_{z_{k}} - \underbrace{\left(\overline{\left(B^{2}\right)^{-1} r_{k}}\right)^{T} A p_{k-1}}_{p_{k-1}^{T} A p_{k-1}} p_{k-1}, \qquad (1.10)$$

and defining the initial iterate as  $u^0 = B^{-1}\tilde{u}^0$  and the initial search direction as  $p_0 = (B^2)^{-1}r_0$ . With this modification, the residuals  $r_k, \tilde{r}_k$ , search directions  $p_k, \tilde{p}_k$ , and iterates  $u^k, \tilde{u}^k$ , of the modified CG algorithm and the preconditioned CG algorithm, verify the relations:  $r_k = B\tilde{r}_k, p_k = B^{-1}\tilde{p}_k, u^k = B^{-1}\tilde{u}^k$ . Indeed, using the assumption that  $r_{k-1} = B\tilde{r}_{k-1}, p_{k-1} = B^{-1}\tilde{p}_{k-1}$ , and  $u^{k-1} = B^{-1}\tilde{u}^{k-1}$  (these assumptions are trivially verifiable at k = 1), we obtain that:

$$\tilde{\alpha}_{k-1} = \frac{\tilde{r}_{k-1}^T \tilde{p}_{k-1}}{\left\| \tilde{p}_{k-1} \right\|_{\tilde{A}}^2} = \frac{r_{k-1}^T B^{-1} B p_{k-1}}{\left( B p_{k-1} \right)^T B^{-1} A B^{-1} B p_{k-1}} = \frac{r_{k-1}^T p_{k-1}}{\left\| p_{k-1} \right\|_{A}^2} = \alpha_{k-1},$$

and subsequently that:  $r_k = B\tilde{r}_k$  and  $u^k = B^{-1}\tilde{u}^k$ . For example,

$$\tilde{u}^k = \tilde{u}^{k-1} + \tilde{\alpha}_{k-1}\tilde{p}_{k-1} = B(u^{k-1} + \alpha_{k-1}p_{k-1}) = Bu^k.$$

This immediately yields that  $r_k = B\tilde{r}_k$ . Then, if  $p_k$  is computed from (1.10) we obtain, multiplying it by B, that:

$$Bp_{k} = B^{-1}r_{k} - \frac{r_{k}^{T}B^{-1}B^{-1}AB^{-1}Bp_{k-1}}{p_{k-1}^{T}BB^{-1}AB^{-1}Bp_{k-1}}Bp_{k-1},$$

or, using the induction assumptions:

$$Bp_k = \tilde{r}_k - \frac{\langle \tilde{r}_k, \tilde{p}_{k-1} \rangle_{\tilde{A}}}{||\tilde{p}_{k-1}||_{\tilde{A}}^2} \tilde{p}_{k-1} = \tilde{p}_k$$

Thus, using an induction argument we establish that  $r_k = B\tilde{r}_k, p_k = B^{-1}\tilde{p}_k, u^k = B^{-1}\tilde{u}^k$ , and therefore,  $\tilde{\alpha}_k = \alpha_k$ , for all positive integer k, Then, it is clear that the modified CG algorithm does not require the knowledge of  $\tilde{A}$  to proceed. It needs only one extra step for the computation of the new search direction. Instead of using  $r_k$  for its computation, we need to use  $z_k$  that is a solution of the system  $B^2 z_k = r_k$ . The conjugate gradient method with preconditioning is given by: Algorithm Preconditioned CGM:

1.  $k \leftarrow 0$ 2.  $r_0 \leftarrow b - Au^0$ 3.  $p_0 = (B^2)^{-1} r_0$ 4. For  $k = 1, 2, \dots, m$ : 5.  $\alpha_{k-1} \leftarrow r_{k-1}^T p_{k-1} / ||p_{k-1}||_A^2$ 6.  $u^k \leftarrow u^{k-1} + \alpha_{k-1}p_{k-1}$ 7.  $r_k \leftarrow r_{k-1} - \alpha_{k-1}Ap_{k-1}$ 8. if  $||r_k||_2 \ge \tau$  then 9. solve  $B^2 z_k = r^k$ 10. / `

$$p_k \leftarrow z_k - \frac{\langle z_k, p_{k-1} \rangle_A}{||p_{k-1}||_A^2} p_{k-1}$$

- 11. k = k + 1
- 12. Go to (5)
- 13. End if
- 14. End

#### Analysis of the CGM

**Lemma 1.7.7** If  $A \in \mathbb{R}^{n \times n}$  is spd, then for m = 0, 1, 2, ... we have

$$\operatorname{span} \{p_0, p_1, \dots, p_m\} = \operatorname{span} \{r_0, r_1, \dots, r_m\} = \operatorname{span} \{r_0, Ar_0, \dots, A^m r_0\}$$

Proof 1. For m = 0, this is trivial.

2. Suppose that for m = k we have

$$\operatorname{span} \{p_0, p_1, \dots, p_k\} = \operatorname{span} \{r_0, r_1, \dots, r_k\} = \operatorname{span} \{r_0, Ar_0, \dots, A^k r_0\}.$$

3. To complete the proof we should show that the same is true for m = k + 1 i.e.

span 
$$\{p_0, p_1, \dots, p_{k+1}\}$$
 = span  $\{r_0, r_1, \dots, r_{k+1}\}$  = span  $\{r_0, Ar_0, \dots, A^{k+1}r_0\}$ .

We have shown that

$$r_{k+1} = r_k - \alpha_k A p_k$$
$$p_k \in \text{span} \{r_0, A r_0, \dots, A^k r_0\}$$
$$A p_k \in \text{span} \{A r_0, A^2 r_0, \dots, A^{k+1} r_0\}$$

this implies

$$r_{k+1} \in \text{span}\{r_0, Ar_0, \dots, A^{k+1}r_0\}$$

so that

$$\operatorname{span}\{r_0, r_1, \dots, r_{k+1}\} \subset \operatorname{span}\{r_0, Ar_0, \dots, A^{k+1}r_0\}$$
(1.11)

#### 1.7. ITERATIVE METHODS

Now,

$$A^{k}r_{0} \in \operatorname{span} \{p_{0}, p_{1}, \dots, p_{k}\}$$
$$A^{k+1}r_{0} \in \operatorname{span} \{Ap_{0}, Ap_{1}, \dots, Ap_{k}\}$$
$$r_{k+1} = r_{k} - \alpha_{k}Ap_{k}$$

So that

$$A^{k+1}r_0 \in \operatorname{span}\left\{\underbrace{Ap_0}_{(r_1-r_0)/\alpha_0}, \underbrace{Ap_1}_{(r_2-r_1)/\alpha_1}, \ldots, \underbrace{Ap_{k-1}}_{(r_k-r_{k-1})/\alpha_{k-1}}, r_{k+1} - r_k\right\} \subset \operatorname{span}\left\{r_0, r_1, \ldots, r_{k+1}\right\}.$$

Therefore

$$A^{k+1}r_0 \in \text{span}\{r_0, r_1, \dots, r_{k+1}\}$$

and by induction

$$\operatorname{span}\left\{r_0, Ar_0, \dots, A^{k+1}r_0\right\} \subset \operatorname{span}\left\{r_0, \dots, r_{k+1}\right\}$$

The last inclusion together with (1.11) prove that:

span  $\{r_0, Ar_0, \dots, A^m r_0\}$  = span  $\{r_0, \dots, r_m\}$ ,  $\forall m \ge 0$ .

$$\operatorname{span} \{p_0, p_1, \dots, p_m\} = \operatorname{span} \{r_0, \dots, r_m\}$$

is proven using

$$p_m = r_m + \beta_{m-1} p_{m-1}$$

and similar arguments as above.

**Definition**  $K_m := \text{span}\{r_0, r_1, \dots, r_{m-1}\}$  is the *m*-th Krylov space of A.

Theorem 1.7.8 If A is spd then:

(A)  $\langle p_k, p_m \rangle_A = 0$  for  $m \neq k$  and (B)  $r_k^T r_m = 0$  for  $m \neq k$ .

**Proof** 1. For  $k, m \leq 1$ :

$$r_1^T r_0 = 0$$
$$\langle p_1, p_0 \rangle_A = 0$$

2. Assume:

(a)  $\langle p_k, p_m \rangle_A = 0$  for  $k \neq m$  and  $k, m \leq l$ 

(b) 
$$r_k^T r_m = 0$$
 for  $k \neq m$  and  $k, m \leq l$ 

3. For m < l:

$$r_l^T p_m = r_l^T (c_0 r_0 + c_2 r_2 + \dots + c_m r_m) = 0$$

from (b) and therefore

$$r_{l+1}^T p_m = \left(r_l - \alpha_l A p_l\right)^T p_m = 0.$$

For m = l we have

$$r_{l+1}^T p_l = \left(r_l - \alpha_l A p_l\right)^T p_l = r_l^T p_l - \alpha_l \left\langle p_l, p_l \right\rangle_A = 0$$

by the definition of  $\alpha_l$ . But  $r_m \in \text{span} \{p_0, \ldots, p_m\}$  for  $m = 1, \ldots, l$  and therefore

 $r_{l+1}^T r_m = 0, \quad m = 0, 1, \dots, l,$ 

which proves (B).

To prove (A), we use:

$$p_{l+1} = r_{l+1} + \beta_l p_l$$
  
$$\langle p_{l+1}, p_m \rangle_A = \langle r_{l+1}, p_m \rangle_A + \beta_l \langle p_l, p_m \rangle_A; \langle p_l, p_m \rangle_A = 0, \quad m = 0, \dots, l-1.$$

 $\operatorname{But}$ 

$$Ap_m \in \operatorname{span}\left\{r_0, r_1, \dots, r_{m+1}\right\}$$

and from (B),

$$\langle r_{l+1}, p_m \rangle_A = 0 \quad \Rightarrow \quad \langle p_{l+1}, p_m \rangle_A = 0 \text{ for } m = 0, 1, \dots, l-1.$$

By the choice of  $p_{l+1}$  in the CGM,

 $\langle p_{l+1}, p_l \rangle_A = 0$ 

and so (A) holds.

From this theorem we can conclude that

 $\dim K_n \le n.$ 

**Corollary 1.7.9** If  $A \in \mathbb{R}^{n \times n}$  is spd then for some  $m \leq n, r_m = 0$ .

**Theorem 1.7.10** If A is spd then  $u^{k+1}$  minimizes the error  $u - u^{k+1} = \epsilon_{k+1}$  over  $K_{k+1}$  in  $|| \cdot ||_A$  i.e.

 $||\epsilon_{k+1}||_A = \min_{v \in K_{k+1}} ||u - v||_A \quad (u : Au = b)$ 

**Proof** Corollary 1.7.9 concluded that there exist some  $m + 1 \le n$  s.t.  $r_{m+1} = 0$ . Without loss of generality we can assume that  $u^0 = 0$  since any change of the initial iterate can be interpreted as a change in the right hand side vector b:  $A(u - u^0) = b - Au^0$ , and therefore it does not affect the estimates that follow. Then the exact solution must be a linear combination of the search directions  $p_0, \ldots, p_m$  i.e.  $u = \sum_{i=0}^m \alpha_i p_i$ . Consider some

 $v \in K_{k+1}, k < m+1$ . It must have the form  $v = \sum_{i=0}^{k} \gamma_i p_i$ . Then we have:

$$||u - v||_A^2 = ||\sum_{i=0}^k (\alpha_i - \gamma_i)p_i + \sum_{i=k+1}^m \alpha_i p_i||_A^2 = \sum_{i=0}^k |\alpha_i - \gamma_i|^2 ||p_i||_A^2 + \sum_{i=k+1}^m |\alpha_i|^2 ||p_i||_A^2,$$

since  $\langle p_i, p_j \rangle_A = 0$  if  $i \neq j$ . Therefore,  $||u - v||_A$  has its minimum for  $\alpha_i = \gamma_i, i = 1, ..., k$ . But  $u^{k+1} = \sum_{i=0}^k \alpha_i p_i$  which concludes the proof.

**Definition**  $\Pi_l$  denotes the set of all polynomials up to order *l*:

$$\Pi_l = \left\{ p^l : a_l x^l + \dots + a_0 \right\}$$

**Corollary 1.7.11** If  $A \in \mathbb{R}^{n \times n}$  is symmetric positive definite with a spectrum  $\{\lambda_i\}_{i=1}^n$  then:

$$\left\|\epsilon_{k+1}\right\|_{A} \le \left\|\epsilon_{0}\right\|_{A} \max_{1 \le j \le n} \left|q\left(\lambda_{j}\right)\right|$$

for all polynomials  $q \in \Pi_{k+1}$  such that q(0) = 1

#### 1.7. ITERATIVE METHODS

**Proof** From theorem 1.7.10 we have that:

$$||\epsilon_{k+1}||_A \le ||u-v||_A, \forall v \in K_{k+1}.$$

But

$$\nu \in K_{k+1} = \operatorname{span}\left\{r_0, Ar_0, \dots, A^k r_0\right\}$$

Without loss of generality, we can assume that  $r_0 = b$  i.e.  $u^0 = 0$  and then

$$v = \alpha_0 b + \alpha_1 A b + \dots + \alpha_k A^k b = p(A) b$$

Therefore

$$||\epsilon_{k+1}||_{A} \le ||u - p(A) b||_{A}, \forall p(A) \in \Pi_{k}$$

but  $u^0 = 0$ , b = Au, and  $\epsilon_0 = u$  so that

$$\begin{aligned} ||\epsilon_{k+1}||_{A} &\leq \left| \left| \epsilon_{0} - p\left(A\right)A\left(u - u^{0}\right) \right| \right|_{A} = \left| \left| \epsilon_{0}\left(I - p\left(A\right)A\right) \right| \right|_{A} \\ &\leq \left| \left| \epsilon_{0} \right| \right| \left| \left| \underbrace{I - p\left(A\right)A}_{\in \Pi_{k+1}} \right| \right|_{A} = \left| \left| \epsilon_{0} \right| \left| \left| \left| q\left(A\right) \right| \right|_{A}, \forall q\left(A\right) \in \Pi_{k+1} \text{ s.t. } q\left(0\right) = 1. \end{aligned}$$

Now we will prove that

$$\left\|\left|q\left(A\right)\right\|\right|_{A} = \max_{j} \left|q\left(\lambda_{j}\right)\right|.$$

Since A is symmetric, there exists an orthonormal set of eigenvectors  $\{v_j\}_{j=1}^n$  of A i.e.

$$v_i^T v_j = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases}$$

Take  $x \in \mathbb{R}^n$  then

$$x = \sum_{i=1}^{n} c_i v_i$$

and therefore  $\left|\left|x\right|\right|_{A}^{2}=\sum\limits_{i=1}^{n}c_{i}^{2}\lambda_{i}$  . Then we have

$$\begin{aligned} ||q(A)||_{A}^{2} &= \sup_{||x||_{A}=1} ||q(A)x||_{A}^{2} = \sup_{\sum_{i=1}^{n} c_{i}^{2}\lambda_{i}=1} \left( q(A)\sum_{j=1}^{n} c_{j}v_{j} \right)^{T} A\left( q(A)\sum_{i=1}^{n} c_{i}v_{i} \right) = \\ &\sup_{\sum_{i=1}^{n} c_{i}^{2}\lambda_{i}=1} \sum_{i=1}^{n} q^{2}\left(\lambda_{i}\right) c_{i}^{2}\lambda_{i} \le \max_{1\le i\le n} q^{2}\left(\lambda_{i}\right). \end{aligned}$$

If  $\max_{1 \le i \le n} q^2(\lambda_i)$  is achieved for some index l then it is clear that the equality in the last relation would be achieved if we take  $x = v_l/\sqrt{\lambda_l}$ , and this concludes the proof of the corollary.

From corollary 1.7.11 it is clear that the sharpest error estimate would be obtained if we pick q to be such that  $\max_{1 \le j \le n} |q(\lambda_j)|$  is minimal with respect to q. Let us denote the minimum and maximum eigenvalues of A by  $\lambda_{min}$  and  $\lambda_{max}$  correspondingly. Before we construct such a polynomial on the interval  $[\lambda_{min}, \lambda_{max}]$  we first consider polynomials on [-1, 1] whose maximum is minimal among all possible polynomials of the same degree i.e. polynomials with a minimal infinity norm. To remind you:

**Definition** The *infinity norm* of a polynomial on [-1, 1] is defined as

$$||p||_{L^{\infty}[-1,1]} = \sup_{z \in [-1,1]} |p(z)|$$

If p(0) = 1 then  $||p||_{L^{\infty}[-1,1]} \ge 1$  so the minimum for such polynomials would be 1.

The polynomials that have a unit infinity norm on [-1, 1] are the *Chebyshev* polynomials. They are defined recursively as follows:

#### Definition

$$T_{0}(z) = 1$$
  

$$T_{1}(z) = z$$
  

$$T_{n+1}(z) = 2zT_{n}(z) - T_{n-1}(z)$$

Alternatively, we may represent the Chebyschev polynomials as

$$T_{k}(z) = \frac{1}{2} \left[ \left( z + \sqrt{z^{2} - 1} \right)^{k} + \left( z - \sqrt{z^{2} - 1} \right)^{k} \right],$$

or

$$T_k(z) = \cos(k \arccos(z)).$$

The polynomials with a minimum infinity norm on  $[\lambda_{min}, \lambda_{max}]$ , such that they equal 1 at 0, are given by the *Shifted Chebyshev* polynomials:

#### Definition

$$\tilde{T}_{k}\left(z\right) = \frac{T_{k}\left(1 - 2\frac{z - \lambda_{\min}}{\lambda_{\max} - \lambda_{\min}}\right)}{T_{k}\left(1 + 2\frac{\lambda_{\min}}{\lambda_{\max} - \lambda_{\min}}\right)}$$

Lemma 1.7.12

$$\left\| \left| \tilde{T}_k\left(z\right) \right\|_{L^{\infty}[\lambda_{\min},\lambda_{\max}]} = \min_{p \in \Pi_k, p(0) = 1} \left\| \left| p(z) \right\|_{L^{\infty}[\lambda_{\min},\lambda_{\max}]}.$$

**Proof** Since  $T_k(z) = \cos(k \arccos(z))$  it is clear that  $T_k$  has k+1 extrema  $T_k(z_i) = (-1)^i$ ,  $i = 0, 1, \ldots, k$  in the nodes  $z_i = \cos(i\pi/k)$  i.e it has k+1 alternating minima and maxima on [-1, 1].  $\tilde{T}_k$  is just a scaled and shifted version of  $T_k$  so it must have exactly the same number of alternating maxima and minima on  $[\lambda_{min}, \lambda_{max}]$ .

Now assume towards contradiction, that there is some  $p_k(z)$  such that  $p_k(0) = 1$  and  $||p_k(z)||_{L^{\infty}[\lambda_{min},\lambda_{max}]} < ||\tilde{T}_k(z)||_{L^{\infty}[\lambda_{min},\lambda_{max}]}$ . Then,  $p_k - \tilde{T}_k$  must also have at least the same number of alternating extrema as  $\tilde{T}_k$  on  $[\lambda_{min}, \lambda_{max}]$  since  $\max |p_k| < \max |\tilde{T}_k|$  on this interval. This means that  $p_k - \tilde{T}_k$  has at least k zeros on  $[\lambda_{min}, \lambda_{max}]$ . But  $p_k(0) - \tilde{T}_k(0) = 0$  and  $0 < \lambda_{min}$ . So we reach the contradiction that the polynomial of k-th degree  $p_k - \tilde{T}_k$  has at least k + 1 zeros which concludes the proof.

From the corollary we have that

$$\left\|\epsilon_{k}\right\|_{A} \leq \left\|\epsilon_{0}\right\|_{A} \max_{\lambda_{i}} \left|p\left(\lambda_{i}\right)\right|$$

for any  $p(x) \in \Pi_k$  such that p(0) = 1. Therefore,

$$\left\| \epsilon_{k} \right\|_{A} \leq \left\| \epsilon_{0} \right\|_{A} \sup_{\lambda_{\min} < z < \lambda_{\max}} \left| p\left( z \right) \right| = \left\| \epsilon_{0} \right\|_{A} \left\| p \right\|_{L^{\infty}[\lambda_{\min}, \lambda_{\max}]}$$

Since p can be any polynomial in  $\Pi_k$  such that p(0) = 1, the sharpest estimate would be obtained for polynomials with a minimum infinity norm i.e. shifted Chebishev polynomials:

$$\left|\left|\epsilon_{k}\right|\right|_{A} \leq \left|\left|\epsilon_{0}\right|\right|_{A} \max_{\lambda_{\min} < z < \lambda_{\max}} \left|\tilde{T}_{k}\left(z\right)\right|$$

#### 1.7. ITERATIVE METHODS

Then,

$$\begin{split} \max_{\lambda_{\min} < z < \lambda_{\max}} \left| \tilde{T}_k(z) \right| &= \max \left| \frac{T_k \left( 1 - 2 \frac{z - \lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} \right)}{T_k \left( 1 + 2 \frac{\lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} \right)} \right| \\ &\leq \frac{1}{\left| T_k \left( 1 + 2 \frac{\lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} \right) \right|} \leq 2 \left( \frac{\sqrt{\lambda_{\max}} - \sqrt{\lambda_{\min}}}{\sqrt{\lambda_{\max}} + \sqrt{\lambda_{\min}}} \right)^k \end{split}$$

where we have used

$$T_k \left( \frac{\lambda_{\max} + \lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} \right) \bigg| = \frac{1}{2} \left( \left( \frac{\sqrt{\lambda_{\max}} + \sqrt{\lambda_{\min}}}{\sqrt{\lambda_{\max}} - \sqrt{\lambda_{\min}}} \right)^k + \left( \frac{\sqrt{\lambda_{\max}} - \sqrt{\lambda_{\min}}}{\sqrt{\lambda_{\max}} + \sqrt{\lambda_{\min}}} \right)^k \right)$$
$$\geq \frac{1}{2} \left( \frac{\sqrt{\lambda_{\max}} + \sqrt{\lambda_{\min}}}{\sqrt{\lambda_{\max}} - \sqrt{\lambda_{\min}}} \right)^k.$$

Then

$$\left|\epsilon_{k}\right|_{A} \leq 2\left|\left|\epsilon_{0}\right|\right|_{A} \left(\frac{\sqrt{\lambda_{\max}} - \sqrt{\lambda_{\min}}}{\sqrt{\lambda_{\max}} + \sqrt{\lambda_{\min}}}\right)^{k}$$

If A is spd then

$$c_{2}(A) = ||A||_{2} ||A^{-1}||_{2} = \frac{\lambda_{\max}}{\lambda_{\min}}$$

and

$$\|\epsilon_k\|_A \le 2 \|\epsilon_0\|_A \left(\frac{\sqrt{c_2(A)} - 1}{\sqrt{c_2(A)} + 1}\right)^k$$

Let  $p(\gamma)$  be the smallest integer k > 0 such that  $||\epsilon_k||_A \leq \gamma ||\epsilon_0||_A$ . Then, from the estimate above (which is sharp) it follows that

$$2\left(\frac{1-z}{1+z}\right)^k \le \gamma, \forall z \in [0,1), \text{ where } z = 1/\sqrt{c_2(A)},$$

or

$$\frac{2}{\gamma} \le \left(\frac{1+z}{1-z}\right)^k.$$

Note that z = 0 corresponds to the limit  $c_2(A) \to \infty$ . Then we get

$$\ln\frac{2}{\gamma} \le k\ln\frac{1+z}{1-z} = 2k(z+\frac{1}{3}z^3+\frac{1}{5}z^5+\dots), \quad z \in [0,1)$$

or

$$\frac{1}{2z}\ln\frac{2}{\gamma} \le k(1 + \frac{1}{3}z^2 + \frac{1}{5}z^4 + \dots) = k\frac{arctanh(z)}{z}.$$

But  $\operatorname{arctanh}(z)/z$  is a monotonically increasing function on [0, 1) and its minimum is 1. Therefore, we have;

$$\frac{1}{2}\sqrt{c_2(A)}\ln\frac{2}{\gamma} \le k,$$

which means that  $p(\gamma) \leq 1/2\sqrt{c_2(A)}\ln(2/\gamma) + 1$ , i.e. for all practical purposes, we need of order of  $\sqrt{c_2(A)}$  iterations for satisfying the convergence condition  $||\epsilon_k||_A \leq \gamma ||\epsilon_0||_A$ .

# CHAPTER 1. SYSTEMS OF LINEAR EQUATIONS

# Chapter 2

# Solutions to Partial Differential Equations

# 2.1 Classification of Partial Differential Equations

Suppose we have a pde of the form

$$\phi\left(x, y, u, \frac{\partial u}{\partial x}, \frac{\partial u}{\partial y}, \frac{\partial^2 u}{\partial x^2}, \frac{\partial^2 u}{\partial x \partial y}, \frac{\partial^2 u}{\partial y^2}\right) = 0$$

then we classify the equation in terms of the second order derivatives, if they are not present we then classify in terms of the first order derivatives.

#### 2.1.1 First Order linear PDEs

$$\alpha\left(t,x\right)\frac{\partial u}{\partial t} + \beta\left(t,x\right)\frac{\partial u}{\partial x} = \gamma\left(t,x\right)$$

Let us try to reduce it to an ODE over some path (t(s), x(s)) in the domain of the equation. We have that:

$$\frac{du}{ds} = \frac{\partial u}{\partial t}\frac{dt}{ds} + \frac{\partial u}{\partial x}\frac{dx}{ds}$$

If we can choose a path  $\alpha = \frac{dt}{ds}$  and  $\beta = \frac{\partial x}{\partial s}$  then

$$\frac{du}{ds} = \gamma$$

These paths are called *characteristics*.

#### 2.1.2 Second Order PDE

$$\alpha \frac{\partial^2 u}{\partial x^2} + \beta \frac{\partial^2 u}{\partial x \partial y} + \gamma \frac{\partial^2 u}{\partial y^2} = \psi \left( x, y, u, \frac{\partial u}{\partial x}, \frac{\partial u}{\partial y} \right)$$

We use a similar idea as in the first order case but now we try to reduce the equation to a system of ODEs for the first derivatives of the solution i.e.

$$\frac{d}{ds}\left(\frac{\partial u}{\partial x}\right) = \frac{\partial^2 u}{\partial x^2}\frac{dx}{ds} + \frac{\partial^2 u}{\partial x \partial y}\frac{dy}{ds}$$
$$\frac{d}{ds}\left(\frac{\partial u}{\partial y}\right) = \frac{\partial^2 u}{\partial x \partial y}\frac{dx}{ds} + \frac{\partial^2 u}{\partial y^2}\frac{dy}{ds}$$
$$\frac{23}{23}$$

Then we may write this system as

$$\begin{bmatrix} \alpha & \beta & \gamma \\ \frac{dx}{ds} & \frac{dy}{ds} & 0 \\ 0 & \frac{dx}{ds} & \frac{dy}{ds} \end{bmatrix} \begin{bmatrix} \frac{\partial^2 u}{\partial x^2} \\ \frac{\partial^2 u}{\partial x \partial y} \\ \frac{\partial^2 u}{\partial y^2} \end{bmatrix} = \begin{bmatrix} \psi \\ \frac{d}{ds} \left( \frac{\partial u}{\partial x} \right) \\ \frac{d}{ds} \left( \frac{\partial u}{\partial y} \right) \end{bmatrix}$$

If we require, similarly to the first-order case, that the original equation is a linear combination of the two equations for the first partial derivative, i. e. this system is linearly dependent, then

$$\alpha \left(\frac{dy}{ds}\right)^2 - \beta \frac{dx}{ds} \frac{dy}{ds} + \gamma \left(\frac{dx}{ds}\right)^2 = 0 \quad \text{and} \quad \alpha \left(\frac{dy}{dx}\right)^2 - \beta \frac{dy}{dx} + \gamma = 0.$$

We classify the equations based on the discriminant

# 2.2 Difference Operators

Suppose that we have a grid of nodes in an interval  $[a, b], \Delta = \{a = x_0 < x_1 < \cdots < x_{N-1} < x_N = b\}$ . The following mapping:  $u_h : \Delta \to \mathbb{R}$  is called a grid function:  $u_h = (u_{h,0}, \ldots, u_{h,N})^T$  with  $u_{h,j} = u_h(x_j)$ . A classical function  $u : [a, b] \to \mathbb{R}$  gives rise to an associated grid function and, abusing notation somewhat, we denote this grid function by u and its value at a given node  $x_j$  by  $u_j$ . In addition we will make use of the following operators on such functions, that are used to approximate the corresponding derivatives:

$$\begin{split} \delta^+ u_j &= \frac{1}{h_{j+1}} \left( u_{j+1} - u_j \right), \qquad h_{j+1} = x_{j+1} - x_j \\ \delta^- u_j &= \frac{1}{h_j} \left( u_j - u_{j-1} \right) \\ \delta u_{j+\frac{1}{2}} &= \frac{1}{h_{j+1}} \left( u_{j+1} - u_j \right), \\ \delta^2 u_j &= \frac{1}{h_{j+\frac{1}{2}}} \left( \frac{u_{j+1} - u_j}{h_{j+1}} - \frac{u_j - u_{j-1}}{h_j} \right), \qquad h_{j+\frac{1}{2}} = \frac{1}{2} \left( h_j + h_{j+1} \right). \end{split}$$

We will also make use of the averaging operator:

$$u_{j+\frac{1}{2}} = \frac{1}{2} \left( u_j + u_{j+1} \right).$$

Assuming that  $h_j = h, \forall j$ , and enough regularity of the function u, and using Taylor expansions, we can derive the following estimates for the error of these approximations:

$$\delta^{+}u_{j} = \frac{1}{h} \left( u_{j+1} - u_{j} \right) = \frac{1}{h} \left( u_{j} + \frac{\partial u_{j}}{\partial x} h + \frac{\partial^{2} u_{j}}{\partial x^{2}} \frac{h^{2}}{2} + \dots - u_{j} \right)$$
$$= \frac{\partial u_{j}}{\partial x} + h \frac{\partial^{2} u_{j}}{\partial x^{2}} + \dots = \frac{\partial u_{j}}{\partial x} + \mathcal{O} \left( h \right)$$
$$\delta^{-}u_{j} = \frac{\partial u_{j}}{\partial x} + \mathcal{O} \left( h \right)$$
$$\delta^{2}u_{j} = \frac{u_{j+1} - 2u_{j} + u_{j-1}}{h^{2}} = \frac{\partial^{2} u_{j}}{\partial x^{2}} + \mathcal{O} \left( h^{2} \right)$$
$$u_{j+\frac{1}{2}} = \frac{1}{2} \left( u_{j} + u_{j+1} \right) = u \left( x_{j+\frac{1}{2}} \right) + \mathcal{O} \left( h^{2} \right).$$

Similar estimates can be derived for non-constant grid size  $h_j$ , however, in the rest of the notes we will consider only equidistant grids. The non-equidistant case requires more careful examination.

#### 2.2. DIFFERENCE OPERATORS

#### 2.2.1 Poisson Equation

$$-\nabla^2 u\left(x,y\right) = f\left(x,y\right), \quad (x,y) \in (0,1) \times (0,1) = \Omega$$
$$u\left(x,y\right) = g\left(x,y\right), \quad (x,y) \in \partial\Omega$$

Consider the grid:

$$\begin{split} \bar{\Omega}_h &= \{ (x_j, y_l) ; x_j = jh, \ y_l = lh, \quad j, l \in 0, \dots, N \} \\ \Omega_h &= \{ (x_j, y_l) ; x_j = jh, \ y_l = lh, 1 \le j, l \le N-1 \} \\ \partial \Omega_h &= \bar{\Omega}_h \setminus \Omega_h, \end{split}$$

where h = 1/N. A second order scheme for an approximation to the solution  $u_h$  is given by:

$$-\left(\delta_x^2 u_{h,j,l} + \delta_y^2 u_{h,j,l}\right) = f_{h,j,l}, \quad (x_j, y_l) \in \Omega_h$$
$$u_{j,l} = g_{h,j,l}, \quad (x_j, y_l) \in \partial\Omega_h$$

where  $f_{h,j,l} = f_{j,l}, (x_j, y_l) \in \Omega_h; g_{h,j,l} = g_{j,l}, (x_j, y_l) \in \partial \Omega_h$  So that

$$-(u_{h,j-1,l}+u_{h,j,l-1}-4u_{h,j,l}+u_{h,j,l+1}+u_{h,j+1,l}) = h^2 f_{h,j,l}$$

The left hand side is a discretization of  $\nabla^2 u$  on the following stencil



where

$$T \equiv \begin{bmatrix} -4 & 1 \\ 1 & -4 & 1 \\ & \ddots \\ & 1 & -4 & 1 \\ & & 1 & -4 \end{bmatrix}, \quad u_{h,j} \equiv \begin{bmatrix} u_{h,j,1} \\ u_{h,j,2} \\ \vdots \\ u_{h,j,N-1} \end{bmatrix}, \quad f_j \equiv \begin{bmatrix} f_{j,1} \\ f_{j,2} \\ \vdots \\ f_{j,N-1} \end{bmatrix}, \quad r_j \equiv f_j + \frac{b_j}{h^2}$$
$$b_j \equiv \begin{bmatrix} g_{j,0} \\ 0 \\ \vdots \\ 0 \\ g_{j,N} \end{bmatrix} \text{ for } j = 2, \dots, N-2, \quad b_1 \equiv \begin{bmatrix} g_{0,1} + g_{1,0} \\ g_{0,2} \\ \vdots \\ g_{0,N-2} \\ g_{0,N-1} + g_{1,N} \end{bmatrix}, \quad b_{N-1} \equiv \begin{bmatrix} g_{N,1} + g_{N-1,0} \\ g_{N,2} \\ \vdots \\ g_{N,N-2} \\ g_{N,N-1} + g_{N-1,N} \end{bmatrix}$$

#### 2.2.2 Neumann Boundary Conditions



Suppose that we need to impose a Neumann condition

$$\left. \frac{\partial u}{\partial n} \right|_{ij} = \left. g \right|_{ij}$$

on the left boundary of the domain sketched above. Then, we write the scheme for all nodes on this boundary and discretize the Neumann condition using a central difference scheme, introducing an additional layer of points (corresponding to level -1 in x direction)

$$\frac{u_{h,1,l} - u_{h,-1,l}}{2h} = g_{0,l}, \quad u_{h,-1,l} + u_{h,0,l-1} - 4u_{h,0,l} + u_{h,0,l+1} + u_{h,1,l} = -h^2 f_{0,l}.$$

In order to eliminate the additional layer of points, we combine these results to get

$$u_{h,0,l-1} - 4u_{h,0,l} + u_{h,0,l+1} + 2u_{h,1,l} = -h^2 f_{0,l} + 2hg_{0,l}$$

# 2.3 Consistency and Convergence

Consider the continuous problem

$$Lu = f \quad \text{in } \Omega$$
$$u = g \quad \text{on } \partial \Omega$$

Now, consider the corresponding discrete problem

$$L_h u_h = f_h \quad \text{in } \Omega_h$$
$$u_h = g_h \quad \text{on } \partial \Omega_h.$$

 $f_h, g_h$  are some sort of approximations of f, g and in these notes we assume that  $f_h = f, g_h = g$  in the nodes of the discretization grid.

**Definition** For a given  $\phi \in \mathcal{C}^{\infty}(\Omega)$  and  $\boldsymbol{x}_h \in \Omega_h$ 

$$\tau_h(\boldsymbol{x}_h) \equiv (L - L_h) \phi(\boldsymbol{x}_h)$$

is called a truncation error of  $L_h$ . The scheme  $L_h$  is consistent with L if

$$\lim_{h\to 0}\tau_h\left(\boldsymbol{x}_h\right)=0$$

for all  $\boldsymbol{x}_{h} \in \Omega_{h}$  and  $\phi \in \mathcal{C}^{\infty}(\Omega)$ . If  $\tau_{h}(\boldsymbol{x}_{h}) = \mathcal{O}(h^{p})$ , then  $L_{h}$  is consistent to order p.

#### 2.3. CONSISTENCY AND CONVERGENCE

#### Proposition 2.3.1

$$-\nabla_h^2 = -\delta_x^2 - \delta_y^2$$
$$-\nabla^2 = -\frac{\partial^2}{\partial x^2} - \frac{\partial^2}{\partial y^2}$$

is consistent to order 2 with

in 
$$\Omega_h$$
.

The consistency of a finite difference grid does not guarantee that the solution to the scheme is approximating well the exact solution of the continuous problem. In order to measure how close the solution to the continuous problem is to the solution of the discrete problem we introduce below the notion of convergence. Since the solution of the continuous problem is a function defined everywhere in  $\Omega$  and the one of the discrete problem is known only in the set of nodes in  $\overline{\Omega}_h$  and we need to compare them somehow, we can either extend the definition of  $u_h$  to the whole set of points inside  $\overline{\Omega}$ , or restrict somehow the function u to the nodes in  $\overline{\Omega}_h$ . The first approach requires the use of some interpolant of  $u_h$ , and then we can compare the two solutions in some continuous norm. In the second approach, we need to use some operator P from the space where u belongs, to the vector space of  $u_h$ . Assume that the grid  $\overline{\Omega}_h$  contains the nodes  $\mathbf{x}_k, k = 0, \ldots, K$ . Then Pu is usually identified with the vector of values of u in all nodes  $\mathbf{x}_k$  i.e.  $Pu = (u(\mathbf{x}_0), \ldots, u(\mathbf{x}_K))^T$ . To be more concise, we often denote  $u(\mathbf{x}_k)$  by  $u_k$  and abusing notations somewhat, we often identify u with Pu if it appears under a discrete norm. Below we follow the second approach and in order to quantify the difference between u and  $u_h$ we introduce the notion of convergence as follows:

**Definition**  $u_h$  converges to u in a given norm ||.|| if  $\epsilon_h = Pu - u_h$  satisfies

$$\lim_{h \to 0} ||\epsilon_h|| = 0.$$

If  $||\epsilon_h|| = \mathcal{O}(h^p)$ , then p is the order of convergence.

In the section on finite difference methods we use exclusively the infinity (or maximum) norm:  $||\epsilon_h|| = \max |\epsilon_{h,j}|$ .

Note that consistency does not guarantee that the difference between the solutions of the continuous and discrete problems also tends to zero. It only guarantees that the action of the difference between the differential and discrete operators on a smooth enough function tends to zero.

Let us go back to the boundary value problem for the Poisson equation

$$-\nabla^2 u = f \quad \text{in } \Omega$$
$$u = g \quad \text{in } \partial\Omega \tag{2.1}$$

with the corresponding discrete problem

$$-\nabla_h^2 u_h = f_h \quad \text{in } \Omega_h$$
$$u_h = g_h \quad \text{in } \partial\Omega_h. \tag{2.2}$$

From the elliptic PDEs theory we know the following maximum principle: If  $\nabla^2 u \ge 0$  in  $\Omega$  then u achieves its maximum on  $\partial\Omega$ . Similarly the discrete Laplacian  $\nabla_h^2$  satisfies a Discrete Maximum Principle:

#### Theorem 2.3.2 (Discrete Maximum Principle) If

$$\nabla_h^2 v_{j,l} \ge 0$$

for all  $(x_j, y_l) \in \Omega_h$ , then

$$\max_{(x_j,y_l)\in\Omega_h} v_{j,l} \le \max_{(x_j,y_l)\in\partial\Omega_h} v_{j,l}$$

Proof

$$\nabla_{h}^{2} v_{j,l} = \frac{v_{j+1,l} + v_{j-1,l} + v_{j,l+1} + v_{j,l-1} - 4v_{j,l}}{h^{2}} \geq 0 \quad \Rightarrow \quad v_{j,l} \leq \frac{1}{4} \left( v_{j+1,l} + v_{j-1,l} + v_{j,l+1} + v_{j,l-1} \right) \quad \forall (x_{j}, y_{l}) \in \Omega_{h}$$

Suppose that the maximum of v is achieved in an internal point  $(x_{\mathcal{J}}, y_{\mathcal{L}}) \in \Omega_h$ , then:

 $v_{\mathcal{J},\mathcal{L}} \geq v_{\mathcal{J}-1,\mathcal{L}}, \quad v_{\mathcal{J},\mathcal{L}} \geq v_{\mathcal{J}+1,\mathcal{L}}, \quad v_{\mathcal{J},\mathcal{L}} \geq v_{\mathcal{J},\mathcal{L}-1}, \quad v_{\mathcal{J},\mathcal{L}} \geq v_{\mathcal{J},\mathcal{L}+1}$ 

so that

$$4v_{\mathcal{J},\mathcal{L}} \ge v_{\mathcal{J}-1,\mathcal{L}} + v_{\mathcal{J}+1,\mathcal{L}} + v_{\mathcal{J},\mathcal{L}-1} + v_{\mathcal{J},\mathcal{L}+1} \quad \Rightarrow \quad 4v_{\mathcal{J},\mathcal{L}} = v_{\mathcal{J}-1,\mathcal{L}} + v_{\mathcal{J}+1,\mathcal{L}} + v_{\mathcal{J},\mathcal{L}-1} + v_{\mathcal{J},\mathcal{L}+1} + v_{\mathcal{J},$$

Therefore, v must attain its maximum in all neighbouring points too. Applying this argument repeatedly to the neighbours we eventually will reach a boundary point.

Corollary 2.3.3 The following two results follow from the discrete maximum principle:

1.

$$\nabla_h^2 u_h = 0 \quad in \ \Omega_h$$
$$u_h = 0 \quad on \ \partial \Omega_h$$

has a unique solution  $u_h = 0$  in  $\Omega_h \cup \partial \Omega_h$ .

2. For given  $f_h$  and  $g_h$ ,

$$\nabla_h^2 u_h = f_h \quad in \ \Omega_h$$
$$u_h = g_h \quad on \ \partial\Omega_h$$

has a unique solution.

**Definition** If

 $v: \ \Omega_h \cup \partial \Omega_h \to \mathbb{R}$ 

then

$$\begin{aligned} ||v||_{\Omega} &\equiv \max_{(x_j, y_l) \in \Omega_h} |v_{j,l}| \\ ||v||_{\partial \Omega} &\equiv \max_{(x_j, y_l) \in \partial \Omega_h} |v_{j,l}| \end{aligned}$$

**Lemma 2.3.4** If  $v_{j,l} = 0$  for all  $(x_j, y_l) \in \partial \Omega_h$  then

$$\left|\left|v\right|\right|_{\Omega} \le M_{\Delta} \left|\left|\nabla_{h}^{2}v\right|\right|_{\Omega}$$

for some  $M_{\Delta} > 0$ , that depends on  $\Omega$  but not on h.

#### $\mathbf{Proof}\ \mathrm{Let}$

$$\left|\left|\nabla_{h}^{2} v\right|\right|_{\Omega} = \nu \quad \Rightarrow \quad -\nu \leq \nabla_{h}^{2} v \leq \nu,$$

and consider the function

$$w: w_{j,l} = \frac{1}{4} \left( x_j^2 + y_l^2 \right) \ge 0.$$

Denote its norm on the boundary by  $M_{\Delta}$  i.e.

$$||w||_{\partial\Omega} = M_{\Delta} > 0.$$

#### 2.4. ADVECTION EQUATION

Note that  $\nabla_h^2 w = 1$ . Then

$$\begin{array}{l} \nabla_h^2 \left( v + \nu w \right) \geq 0 \\ \nabla_h^2 \left( v - \nu w \right) \leq 0 \end{array} \right\} \ \, {\rm in} \ \, \Omega \end{array}$$

From the maximum principle and since  $v_h|_{\partial\Omega_h} = 0$  we have:

$$\begin{aligned} v_{j,l} &\leq v_{j,l} + \nu w_{j,l} \leq \nu \left| |w_{j,l}| \right|_{\partial\Omega} = M_{\Delta} \left| \left| \nabla_{h}^{2} v_{j,l} \right| \right|_{\Omega} \\ v_{j,l} &\geq v_{j,l} - \nu w_{j,l} \geq -\nu \left| |w_{j,l}| \right|_{\partial\Omega} = -M_{\Delta} \left| \left| \nabla_{h}^{2} v_{j,l} \right| \right|_{\Omega} \end{aligned}$$

 $\forall (x_j, y_l) \in \Omega_h$ 

**Theorem 2.3.5 (error estimate)** Let u be the solution of the boundary value problem (2.1) and  $u_h$  is a solution to the discrete analog (2.2). Assuming that  $u \in C^4(\overline{\Omega})$ , then there exists a constant k > 0 such that:

$$||u - u_h||_{\Omega} \le kMh^2$$

where

$$M \equiv \max\left\{ \left\| \left| \frac{\partial^4 u}{\partial x^4} \right| \right|_{L^{\infty}(\Omega)}, \left\| \left| \frac{\partial^4 u}{\partial y^4} \right| \right|_{L^{\infty}(\Omega)} \right\}$$

**Proof** From the proposition:

$$\left(\nabla_h^2 - \nabla^2\right) u_{j,l} = \frac{h^2}{12} \left[\frac{\partial^4 u}{\partial x^4} \left(\xi_j, y_l\right) + \frac{\partial^4 u}{\partial y^4} \left(x_j, \eta_l\right)\right]$$

for some  $\xi_j, \eta_l$  s.t.  $x_{j-1} \leq \xi_j \leq x_{j+1}, y_{l-1} \leq \eta_l \leq y_{l+1}$ . Taking into account that u is the exact solution we have:

$$-\nabla_h^2 u_{j,l} = f_{j,l} - \frac{h^2}{12} \left[ \frac{\partial^4 u}{\partial x^4} \left( \xi_j, y_l \right) + \frac{\partial^4 u}{\partial y^4} \left( x_j, \eta_l \right) \right]$$

Subtract  $-\nabla_h^2 u_{h,j,l} = f_{h,j,l} = f_{j,l}$  and note that  $u - u_h = 0$  on  $\partial \Omega_h$ . Then

$$\nabla_{h}^{2}\left(u_{j,l}-u_{h,j,l}\right) = \frac{h^{2}}{12} \left[\frac{\partial^{4}u}{\partial x^{4}}\left(\xi_{j}, y_{l}\right) + \frac{\partial^{4}u}{\partial y^{4}}\left(x_{j}, \eta_{l}\right)\right]$$

From the Lemma we get

$$\left|\left|u-u_{h}\right|\right|_{\Omega} \leq M_{\Delta}\left|\left|\nabla_{h}^{2}\left(u-u_{h}\right)\right|\right|_{\Omega} \leq \underbrace{\frac{2M_{\Delta}}{12}}_{=k}Mh^{2} \quad \blacksquare$$

# 2.4 Advection Equation

Consider the initial value problem

$$\frac{\partial u}{\partial t} + v \frac{\partial u}{\partial x} = 0, \quad -\infty < x < \infty, \quad t > 0$$

with v > 0, u(0, x) = f(x) for  $-\infty < x < \infty$ . The exact solution is given by

$$u\left(t,x\right) = f\left(x - vt\right)$$

Consider a grid  $\Delta \equiv \Delta_t \times \Delta_x$ :

$$\Delta_t \equiv \{nk : n = 0, 1, 2, \ldots\}$$
$$\Delta_x \equiv \{jh : j = 0, \pm 1, \pm 2, \ldots\}$$

One possible discretization is

$$\left(\delta_t^+ + v\delta_x^-\right)u_{h,j}^n = \frac{u_{h,j}^n - u_{h,j}^{n-1}}{k} + v\frac{u_{h,j}^{n-1} - u_{h,j-1}^{n-1}}{h} = 0$$

or equivalently

$$u_{h,j}^{n} = u_{h,j}^{n-1} - C\left(u_{h,j}^{n-1} - u_{h,j-1}^{n-1}\right)$$

This is called the FB (Forwards Backwards) scheme where the Courant number is defined as



Assume for simplicity that v = 1. As suggested by the exact solution u(t,x) = f(x-t), the solution at any point (t,x) depends only on the value of the initial condition at the point where the characteristic (having a slope of 1) crosses the vertical axes (see the figure above). The characteristic curve therefore is called a domain of dependance of the solution. On the other hand, the numerical solution at each point  $(t^n, x_j)$  depends only of the solution at  $(t^{n-1}, x_{j-1})$  and  $(t^{n-1}, x_j)$ . These values themselves depend on values at time  $t^{n-2}$  etc, so that the value at  $(t^n, x_j)$  depends on the values in a right angle triangle whose hypothenuse has a slope of C (see again the same figure), called numerical domain of dependance. Clearly, if the Courant number C is greater than 1, the numerical domain of dependance will not contain the point where the characteristic crosses the x-axes i.e. it will not account for the initial data that determines the exact solution at a given point. This should create troubles with the numerical solution and these troubles are quantified by the notion of stability. Informally speaking, the solution is stable if the numerical domain of dependance contains the exact domain of dependance. We will formally illustrate these concepts below.

Define S (shift operator) such that

$$Su_{h,j}^n = u_{h,j+1}^n$$

$$u_{h,j}^{n} = (1 - C + CS^{-1}) u_{h,j}^{n-1} = (1 - C + CS^{-1})^{2} u_{h,j}^{n-2} = \dots = (1 - C + CS^{-1})^{n} u_{h,j}^{0} = (1 - C + CS^{-1})^{n} f_{j}$$
$$= \sum_{m=0}^{n} \binom{n}{m} (1 - C)^{m} (CS^{-1})^{n-m} f_{j} = \sum_{m=0}^{n} \binom{n}{m} (1 - C)^{m} C^{n-m} f_{j-n+m}$$

Now, suppose that the initial condition is perturbed with an error  $\epsilon_j$  i.e. instead of  $f_j$  the initial condition is given by  $\hat{f}_j = f_j + \epsilon_j$ . Then, the perturbed solution  $\hat{u}_h$  is given by:

$$\hat{u}_{h,j} = \sum_{m=0}^{n} \binom{n}{m} (1-C)^m C^{n-m} (f_{j-n+m} + \epsilon_{j-n+m})$$

#### 2.5. VON NEUMANN STABILITY ANALYSIS

and we have

$$\begin{aligned} \left| u_{h,j}^{n} - \hat{u}_{h,j}^{n} \right| &= \left| \sum_{m=0}^{n} \binom{n}{m} \left( 1 - C \right)^{m} C^{n-m} \epsilon_{j-n+m} \right| \\ &\leq \sum_{m=0}^{n} \binom{n}{m} \left| 1 - C \right|^{m} C^{n-m} \left| \epsilon_{j-n+m} \right| \leq ||\epsilon||_{\infty} \sum_{m=0}^{n} \binom{n}{m} \left| 1 - C \right|^{m} C^{n-m} = ||\epsilon||_{\infty} \left( ||1 - C| + C \right)^{n} \end{aligned}$$

So that

 $|u_{h,j}^{n} - \hat{u}_{h,j}^{n}| \le ||\epsilon||_{\infty} (|1 - C| + C)^{n}$ 

Then, if  $C \leq 1$ :

$$\left|u_{h,j}^{n} - \hat{u}_{h,j}^{n}\right| \le \left|\left|\epsilon\right|\right|_{\infty}$$

However, if C > 1:

$$u_{h,j}^n - \hat{u}_{h,j}^n \Big| \le ||\epsilon||_{\infty} \left(|1 - C| + C\right)^n = ||\epsilon||_{\infty} \left(2C - 1\right)^n \stackrel{n \to \infty}{\to} \infty$$

**Theorem 2.4.1** If  $C \leq 1$  and  $u \in C^2((0,t] \times \mathbb{R})$  then the FB scheme is convergent. That is

$$E^n = \left\| u^n - u^n_h \right\|_{\infty} \le t^n \mathcal{O}\left(k + h\right)$$

where  $t^n = kn$ .

Proof

$$\begin{split} \delta_t^+ u\left(t^n, x_j\right) &= \frac{\partial u}{\partial t}\left(t^n, x_j\right) + \frac{k}{2} \frac{\partial^2 u}{\partial t^2}\left(\eta, x_j\right) \\ \delta_x^- u\left(t^n, x_j\right) &= \frac{\partial u}{\partial x}\left(t^n, x_j\right) - \frac{h}{2} \frac{\partial^2 u}{\partial x^2}\left(t^n, \xi\right) \\ \delta_t^+ u + v \delta_x^- u &= -\tau_{k,h}\left(t^n, x_j\right) = \mathcal{O}\left(k+h\right) \\ \delta_t^+ u_{h,j}^h + v \delta_x^- u_{h,j}^n &= 0 \end{split}$$

and letting  $\epsilon_h = u - u_h$  we get the following estimate:

$$\left|\epsilon_{h,j}^{n+1}\right| \leq \left|\left(1-C\right)\epsilon_{h,j}^{n}\right| + C\left|\epsilon_{h,j-1}^{n}\right| + k\left|\tau_{k,h}\left(t^{n},x_{j}\right)\right|$$

Let

$$E^{n} = \left| \left| \epsilon_{h}^{n} \right| \right|_{\infty}, \quad T^{n} = \max_{j} \left| \tau_{k,h} \left( t^{n}, x_{j} \right) \right|$$

Then from the last inequality we can conclude that:

$$E^{n+1} \le (1-C) E^n + CE^n + k |\tau_{k,h}(t^n, x_j)| \le E^n + kT^n \le E^{n-1} + kT^n + kT^{n-1} \le E^0 + k \sum_{m=1}^n T^m + kT^m + kT^{n-1} \le E^0 + k \sum_{m=1}^n T^m + kT^m + k$$

If  $E^0 = 0$  then

$$E^{n+1} \le k \sum_{m=1}^{n} T^m \le nk \max_{1 \le m \le n} T^m = t^n \max_m T^m = t^n \mathcal{O}(k+h)$$

# 2.5 Von Neumann Stability Analysis

**Definition** Discrete  $L^2$  norm:

$$||v||_2^2 \equiv h \sum_{j=-\infty}^{\infty} |v_j|^2$$

**Definition** A FD (Finite Difference) scheme for a time dependent PDE is *stable* if for any time T > 0, there exists K > 0 such that for any initial data  $u_h^0$ , the sequence  $\{u_h^n\}$  satisfies:

$$\left|\left|u_{h}^{n}\right|\right|_{2} \leq K\left|\left|u_{h}^{0}\right|\right|_{2}$$

for  $0 \le nk \le T$ . The constant K is independent of the time and space mesh size!

Consider the following grid on the real axes

$$\Delta_x = \{\dots, -2h, -h, 0, h, 2h, \dots\}$$

Then the Fourier Transform of a given grid function v is defined as

$$(Fv)(\tau) \equiv \frac{h}{\sqrt{2\pi}} \sum_{j=-\infty}^{\infty} e^{-ijh\tau} v_j$$

Inverse transform:

$$v_j = \frac{1}{\sqrt{2\pi}} \int_{-\frac{\pi}{h}}^{\frac{\pi}{h}} e^{ijh\xi} \left(Fv\right)\left(\xi\right) d\xi$$

Theorem 2.5.1 (Parseval Identity)

$$||Fv||_{L^{2}\left[-\frac{\pi}{h},\frac{\pi}{h}\right]}^{2} := \int_{-\frac{\pi}{h}}^{\frac{\pi}{h}} (Fv)^{2} d\xi = ||v||_{2}^{2}$$

**Corollary 2.5.2** A finite difference scheme is stable if and only if  $\exists K > 0$ , independent of k, h, s.t.:

$$\left|\left|Fu_{h}^{n}\right|\right|_{L^{2}\left[-\frac{\pi}{h},\frac{\pi}{h}\right]} \leq K\left|\left|Fu_{h}^{0}\right|\right|_{L^{2}\left[-\frac{\pi}{h},\frac{\pi}{h}\right]}$$

Consider the forward backwards scheme. Since we have:

$$u_{h,j}^{n} = (1-C) u_{h,j}^{n-1} + C u_{h,j-1}^{n-1}$$
$$u_{h,j}^{n} = \frac{1}{\sqrt{2\pi}} \int_{-\frac{\pi}{h}}^{\frac{\pi}{h}} e^{ijh\xi} \left(F u_{h}^{n}\right) \left(\xi\right) d\xi$$
$$u_{h,j-1}^{n-1} = \frac{1}{\sqrt{2\pi}} \int_{-\frac{\pi}{h}}^{\frac{\pi}{h}} e^{i(j-1)h\xi} \left(F u_{h}^{n-1}\right) \left(\xi\right) d\xi,$$

-1

4

we obtain:

$$Fu_{h}^{n} = \underbrace{\left(1 - C + Ce^{-ih\xi}\right)}_{\text{Amplification Factor}} \left(Fu_{h}^{n-1}\right)(\xi)$$

$$A(h\xi) = 1 - C + Ce^{-ih\xi}$$

Then, since we have

$$||Fu_{h}^{n}||_{L^{2}\left[-\frac{\pi}{h},\frac{\pi}{h}\right]}^{2} = \int_{-\frac{\pi}{h}}^{\frac{\pi}{h}} A\left(h\xi\right)^{2} \left(Fu_{h}^{n-1}\right)^{2} d\xi,$$

if  $|A(h\xi)| \leq 1$ , we obtain that:

$$||Fu_h^n||_{L^2\left[-\frac{\pi}{h},\frac{\pi}{h}\right]} \le ||Fu_h^{n-1}||_{L^2\left[-\frac{\pi}{h},\frac{\pi}{h}\right]}$$

and recursive application gives:

$$||Fu_h^n||_{L^2\left[-\frac{\pi}{h},\frac{\pi}{h}\right]} \le \left| |Fu_h^0||_{L^2\left[-\frac{\pi}{h},\frac{\pi}{h}\right]} \right|$$

Since

$$A\left(\theta\right)=1-C+Ce^{-i\theta}$$

$$|A(\theta)|^{2} = (1 - C + C\cos\theta)^{2} + (-C\sin\theta)^{2}$$
  
= 1 - 2C + C<sup>2</sup> + 2C cos  $\theta$  - 2C<sup>2</sup> cos  $\theta$  + C<sup>2</sup> cos<sup>2</sup>  $\theta$  + C<sup>2</sup> sin<sup>2</sup>  $\theta$   
= 1 - 2C (1 - cos  $\theta$ ) + 2C<sup>2</sup> (1 - cos  $\theta$ ) = 1 + 2C (1 - cos  $\theta$ ) (C - 1)  
 $\leq$  1 + 4C (C - 1)

it follows that if C > 1 then  $|A(\theta)| > 1$ , and if  $C \le 1$  then  $|A(\theta)| \le 1$ . Also,

$$||1 - C + Ce^{-i\theta}|| \le |1 - C| + C$$

If a numerical scheme requires a solution of a linear system with a non-diagonal matrix then we call it an *implicit* scheme, otherwise it is called an *explicit* scheme.

Consider the following "Backwards Backwards" scheme

$$\frac{u_{h,j}^{n} - u_{h,j}^{n-1}}{k} + v \frac{u_{h,j}^{n} - u_{h,j-1}^{n}}{h} = 0$$
$$(1+C) u_{h,j}^{n} - Cu_{h,j-1}^{n} = u_{h,j}^{n-1}$$

This is an *implicit* scheme.

In order to study its stability (and the stability of any other FD scheme), it is enough, instead of considering the full Fourier transform of the solution, to consider the "action" of the scheme on a single mode of the form:

 $w^n e^{ijh\xi}$ 

then plugging into the BB scheme:

$$w^{n} \underbrace{\left[1 + C\left(1 - e^{-ih\xi}\right)\right]}_{A^{-1}(h\xi)} = w^{n-1}$$
$$w_{n} = A\left(h\xi\right)w^{n-1}$$
$$A\left(\theta\right) = \left[1 + C\left(1 - e^{-i\theta}\right)\right]^{-1}$$

Then,

$$\left|1 + C\left(1 - e^{-ih\xi}\right)\right| \ge |1 + C| - \left|Ce^{-ih\xi}\right| = 1$$

so that

$$\left|A\left(\theta\right)\right| \leq 1 \quad \forall \theta$$

therefore the BB scheme is unconditionally stable. We may outline a procedure for von Neumann stability analysis:

- 1. Consider a single mode  $w^n e^{ijh\xi}$ .
- 2. Derive conditions for A to satisfy  $|A(h\xi)| \leq 1$

## 2.6 Sufficient Conditions for Convergence

More generally, if we approximate the solution of a well-posed initial-boundary or boundary value problem comprised by the equation:

$$Lu = f \text{ in } \Omega$$

and some initial and/or boundary conditions:

lu = g,

by a numerical scheme with a grid size h (this can be a vector of grid sizes in each direction of discretization):

$$L_h u_h = f_h \text{ in } \Omega_h$$

$$l_h u_h = g_h$$
(2.3)

we need to generalize the notion of stability as follows. Assume that  $u_h$ ,  $f_h$ , and  $g_h$  belong to three Banach spaces  $B_h^1$ ,  $B_h^2$ , and  $B_h^3$ , with norms  $\|.\|_{h,1}$ ,  $\|.\|_{h,2}$ , and  $\|.\|_{h,3}$  correspondingly. Then we have the following generalized definition of stability:

**Definition** (Stable scheme) If the scheme (2.3) has a unique solution  $u_h$ , it is stable if  $\exists C_1, C_2 > 0$ , independent of h, s.t.:

$$||u_h||_{h,1} \le C_1 ||f_h||_{h,2} + C_2 ||g_h||_{h,3}.$$

The essence of this definition is that it guarantees that a small perturbation in the data of the problem i.e. in  $f_h, g_h$  leads to a small perturbation in the solution  $u_h$ .

We can also generalize the notion of consistency as follows. Suppose that  $P^1$  is an operator from a given Banach space  $B^1$ , containing the solution of the continuous problem u, into  $B_h^1$ . Then the discrete operators  $L_h$  and  $l_h$  are consistent with the continuous counterparts if:

$$L_h(P^1u - u_h) = \phi_h \text{ in } \Omega_h$$
  

$$l_h(P^1u - u_h) = \psi_h,$$
(2.4)

and

$$\|\phi_h\|_{h,2} \xrightarrow[|h|\to 0]{} 0, \quad \|\psi_h\|_{h,3} \xrightarrow[|h|\to 0]{} 0.$$

Note that h can be a vector of grid sizes, if the problem involves more than one variable. If  $\phi_h = O(|h|^m), \psi_h = O(|h|^k)$ , then the scheme is consistent to order min(k,m). Note that this definition is consistent with the previous definition that we used, provided that the the data of the discrete problem exactly matches the data of the continuous problem in the points of the grid.

The following theorem states that consistency and stability are sufficient for obtaining convergence.

**Theorem 2.6.1** Given a scheme that is consistent in a norm  $\|.\|_{h,1}$ , then it is convergent in this norm, if it is stable. If the scheme is consistent to order k then it is convergent to order k.

**Proof** If the scheme is consistent we have that:

$$L_h(P^1u - u_h) = \phi_h \text{ in } \Omega_h$$
$$l_h(P^1u - u_h) = \psi_h.$$

The stability guarantees that:

$$||P^{1}u - u_{h}||_{h,1} \le C_{1} ||\phi_{h}||_{h,2} + C_{2} ||\psi_{h}||_{h,3},$$

and this immediately yields convergence in the norm  $\|.\|_{h,1}$ .

It appears that the stability and consistency are also necessary conditions for convergence, as stated in the following theorem due to P. Lax. We will skip the proof in these notes.

**Theorem 2.6.2 (Lax Theorem)** Given a consistent scheme for a well-posed initial boundary value problem, then stability is a necessary and sufficient condition for convergence.

# 2.7 Parabolic PDEs – Heat Equation

$$\begin{cases} \frac{\partial u}{\partial t} = D \frac{\partial^2 u}{\partial x^2}, \quad (t, x) \in (0, T] \times \mathbb{R} \\ u(0, x) = f(x), \quad x \in \mathbb{R} \end{cases}$$

Which has the solution

$$u(t,x) = \frac{1}{\sqrt{4\pi Dt}} \int_{-\infty}^{\infty} f(\xi) e^{-\frac{(x-\xi)^2}{4Dt}} d\xi$$

#### 2.7.1 FC Scheme

This is an explicit scheme.

$$\begin{split} \delta_t^+ u_{h,j}^n &- D \delta_x^2 u_{h,j}^n = 0 \\ u_{h,j}^{n+1} &= (1 - 2\Gamma D) \, u_{h,j}^n + \Gamma D \left( u_{h,j+1}^n + u_{h,j-1}^n \right), \qquad \Gamma = \frac{k}{h^2} \end{split}$$

where  $\Gamma$  is called a grid ratio.



The consistency error is found using

$$\partial_t u - \delta_t^+ u = \mathcal{O}(k)$$
$$\partial_{xx} u - \delta_x^2 u = \mathcal{O}(h^2)$$

so that

$$au_{k,h} = \mathcal{O}\left(k+h^2\right).$$

Stability is a central concern of parabolic and hyperbolic PDEs. Using the Neumann analysis idea we substitute the Fourier mode:

$$v_i^n = w^n e^{ij\theta}$$

into the scheme to get:

$$\frac{1}{k} \left( w^{n+1} - w^n \right) e^{ij\theta} = \frac{D}{h^2} \left( w^n e^{i\theta} - 2w^n + w^n e^{-i\theta} \right) e^{ij\theta}$$

This gives:

$$w^{n+1} = \underbrace{\left(1 - 4D\Gamma\sin^2\frac{\theta}{2}\right)}_{A(\theta) \le 1} w^n$$

So that

$$0 < D\Gamma \leq \frac{1}{2} \quad \Leftrightarrow \quad |A(\theta)| \leq 1 \quad \Rightarrow \quad \text{stability}$$

and we require:

$$\frac{Dk}{h^2} \leq \frac{1}{2} \quad \Rightarrow \quad k \leq \frac{1}{2} \frac{h^2}{D}$$

and we note that k and h are not of the same order. So, this is a conditionally stable scheme, but it is explicit and this is to be expected.

In this problem the domain  $\Omega$  is the entire real axes and therefore it is convenient to use for measuring the error the following  $\infty$ -norm:  $||v_h||_{\Omega} = \max_{-\infty < j < \infty} |v_{h,j}|$ . Then the following theorem provides the convergence estimate for the FC scheme.

**Theorem 2.7.1** If  $D\Gamma \leq \frac{1}{2}$ , the exact solution u is sufficiently smooth, and if  $u_h$  satisfies exactly the initial condition then:

$$||u - u_h||_{\Omega} = ||\epsilon_h||_{\Omega} = \mathcal{O}\left(k + h^2\right)$$

**Proof** From the scheme, if the exact solution u is sufficiently smooth<sup>\*</sup>, we have:

$$u_{j}^{n+1} = D\Gamma u_{j-1}^{n} + (1 - 2D\Gamma) u_{j}^{n} + D\Gamma u_{j+1}^{n} + \mathcal{O}\left(k^{2} + kh^{2}\right)$$

so that:

$$\epsilon_{h,j}^{n+1} = \underbrace{D\Gamma}_{\geq 0} \epsilon_{h,j-1}^{n} + \underbrace{(1-2D\Gamma)}_{\geq 0} \epsilon_{h,j}^{n} + \underbrace{D\Gamma}_{\geq 0} \epsilon_{h,j+1}^{n} + \mathcal{O}\left(k^{2} + kh^{2}\right)$$

by the triangle inequality,

$$\left|\left|\epsilon_{h}^{n+1}\right|\right|_{\Omega} \leq \left|\left|\epsilon_{h}^{n}\right|\right|_{\Omega} + \mathcal{O}\left(k^{2} + kh^{2}\right) \leq \left|\left|\epsilon_{h}^{0}\right|\right|_{\Omega} + (n+1)\mathcal{O}\left(k^{2} + kh^{2}\right) = \underbrace{(n+1)k}_{t^{n+1}}\mathcal{O}\left(k + h^{2}\right)$$

Note that the right hand side of this estimate will blow up as  $t \to \infty$  for fixed k, h.

#### 2.7.2 BC scheme



This is a  $\mathcal{O}(k+h^2)$  consistent scheme. It can be rewritten as:

$$-D\Gamma u_{h,j-1}^{n+1} + (1+2\Gamma D) u_{h,j}^{n+1} - D\Gamma u_{h,j+1}^{n+1} = u_{h,j}^{n}$$

So we must invert a tridiagonal matrix. Now we will study its stability. Substituting

$$v_i^n = w^n e^{ij\ell}$$

into the scheme we have:

$$w^{n+1} - w^n = D\Gamma \left( e^{i\theta} - 2 + e^{-i\theta} \right) w^{n+1}$$

So that:

$$w^{n+1} = \underbrace{\frac{1}{1 + 4D\Gamma \sin^2 \frac{\theta}{2}}}_{A(\theta) \le 1} w^n$$

So, it is clear that the BC scheme is unconditionally stable.

<sup>\*</sup>This term is used to require enough of the derivatives of the exact solution to be bounded, so that the coefficients in the function  $\mathcal{O}\left(k^2 + kh^2\right)$  appearing in the consistency estimate below to be finite.

#### **Crank-Nicolson Scheme** 2.7.3

It is given by

$$\delta_t^{-} u_{h,j}^{n+1} = \frac{1}{2} D\left(\delta_x^2 u_{h,j}^{n+1} + \delta_x^2 u_{h,j}^n\right),$$

or

$$\frac{u_{h,j}^{n+1} - u_{h,j}^{n}}{k} = \frac{D}{2h^2} \left( u_{h,j-1}^{n+1} - 2u_{h,j}^{n+1} + u_{h,j+1}^{n+1} + u_{h,j-1}^{n} - 2u_{h,j}^{n} + u_{h,j+1}^{n} \right)$$

Consistency:

$$\frac{D}{2} \left( \delta_x^2 u_j^{n+1} + \delta_x^2 u_j^n \right) = \frac{D}{2} \left( \frac{\partial^2 u}{\partial x^2} \Big|_j^{n+1} + \mathcal{O} \left( h^2 \right) + \frac{\partial^2 u}{\partial x^2} \Big|_j^n + \mathcal{O} \left( h^2 \right) \right)$$
$$= D \left. \frac{\partial^2 u}{\partial x^2} \Big|_j^{n+\frac{1}{2}} + \mathcal{O} \left( k^2 + h^2 \right)$$
$$\frac{u_j^{n+1} - u_j^n}{k} = \left. \frac{\partial u}{\partial t} \right|_j^{n+\frac{1}{2}} + \mathcal{O} \left( k^2 \right)$$
$$t_{n+1} - \left. - \frac{1}{k} \right|_j^{n+\frac{1}{2}} + \mathcal{O} \left( k^2 \right)$$
$$t_{n-1} - \left. - \frac{1}{k} \right|_j^{n+\frac{1}{2}} + \left. \frac{1}$$

Neumann stability analysis: First, rewrite the scheme in the following two-stage form

$$\frac{u_{h,j}^{n+\frac{1}{2}} - u_{h,j}^{n}}{\frac{k}{2}} = D \frac{u_{h,j-1}^{n} - 2u_{h,j}^{n} + u_{h,j+1}^{n}}{h^{2}} \to A_{1}\left(\theta\right) = 1 - 4D\frac{\Gamma}{2}\sin^{2}\frac{\theta}{2}$$
$$\frac{u_{h,j}^{n+1} - u_{h,j}^{n+\frac{1}{2}}}{\frac{k}{2}} = D \frac{u_{h,j-1}^{n+1} - 2u_{h,j}^{n+1} + u_{h,j+1}^{n+1}}{h^{2}} \to A_{2}\left(\theta\right) = \frac{1}{1 + 4D\frac{\Gamma}{2}\sin^{2}\frac{\theta}{2}}.$$

 $x_{j+1}$ 

Then,

$$w^{n+\frac{1}{2}} = A_1(\theta) w^n \quad w^{n+1} = A_2(\theta) w^{n+\frac{1}{2}} = \underbrace{A_1 A_2}_{A(\theta)} w^n$$

So that  $A(\theta) \leq 1$ , so the system is unconditionally stable, and  $\mathcal{O}(h^2 + k^2)$ .

#### 2.7.4 Leapfrog Scheme



As we use central difference in both directions, the scheme is  $\mathcal{O}(k^2 + h^2)$ . We may write the scheme as:

$$\frac{u_{h,j}^{n+1} - u_{h,j}^{n-1}}{2k} = D \frac{u_{h,j-1}^n - 2u_{h,j}^n + u_{h,j+1}^n}{h^2}$$

Now we preform stability analysis:  $v_j^n = w^n e^{ij\theta}$ 

$$w^{n+1} - w^{n-1} = 2D\Gamma \left( e^{-i\theta} - 2 + e^{i\theta} \right) w^n$$

Then, assume that:

$$w^{n} = A(\theta) w^{n-1}$$
 and  $w^{n+1} = A(\theta) w^{n}$ 

so that

$$(A^2 - 1) w^{n-1} = 4D\Gamma(\cos\theta - 1) Aw^{n-1}$$

this gives

$$A_{1,2} = -4D\Gamma \sin^2 \frac{\theta}{2} \pm \sqrt{1 + 16D^2\Gamma^2 \sin^4 \frac{\theta}{2}}$$

clearly  $|A_2| \ge 1$ , and  $|A_2| > 1$  for some  $\theta$ , so the scheme is unconditionally unstable.

#### 2.7.5 DuFort-Frankel Scheme

$$\delta_t u_h^n = \frac{D}{h^2} \left[ u_{h,j+1}^n - \left( u_{h,j}^{n+1} + u_{h,j}^{n-1} \right) + u_{h,j-1}^n \right]$$

Since

$$\frac{1}{2}\left(u_{h,j}^{n+1} + u_{h,j}^{n-1}\right) = u_{h,j}^{n} + \mathcal{O}\left(k^{2}\right)$$

this scheme can be considered as a stabilized version of the Leapfrog scheme. The truncation error is then:

$$\tau_{k,h}(t,x) = 2\left(\frac{k}{h}\right)^2 \frac{\partial^2 u}{\partial t^2}(t_n,x_j) + \frac{k^3}{3} \frac{\partial^2 u}{\partial t^2}(t_n,x_j) - \frac{h^2}{12} \frac{\partial^2 u}{\partial x^2}(t_n,x_j) + \mathcal{O}\left(k^3 + h^3\right)$$

and the system is conditionally consistant if  $\frac{k}{h} \to 0$ . If k, h go to zero but  $\frac{k}{h} \to c^2$ , then DF is consistent with:

$$\frac{\partial u}{\partial t} + c^2 \frac{\partial^2 u}{\partial t^2} = D \frac{\partial^2 u}{\partial x^2}$$

The stability analysis reveals that this scheme is unconditionally stable.

# 2.8 Advection-Diffusion Equation

$$\frac{\partial u}{\partial t} + v \frac{\partial u}{\partial x} = D \frac{\partial^2 u}{\partial x^2} \qquad (t, x) \in (0, T] \times (0, L)$$

and v, D > 0. Then define  $\tau = \frac{vt}{L}$ ,  $\xi = \frac{x}{L}$ , and Peclet number  $Pe = \frac{vL}{D}$ . So the PDE becomes:

$$\frac{\partial u}{\partial \tau} + \frac{\partial u}{\partial \xi} = \frac{1}{Pe} \frac{\partial^2 u}{\partial \xi^2}, \quad (\tau, \xi) \in \left(0, \frac{vT}{L}\right] \times (0, 1)$$

#### 2.8.1 FC Scheme

$$\frac{u_{h,j}^{n+1} - u_{h,j}^{n}}{k} + \frac{u_{h,j+1}^{n} - u_{h,j-1}^{n}}{2h} - \frac{1}{Pe} \frac{u_{h,j-1}^{n} - 2u_{h,j}^{n} + u_{h,j+1}^{n}}{h^{2}} = 0$$

This scheme is consistant to  $\mathcal{O}(k+h^2)$ . Now preform stability analysis:  $v_i^n = w^n e^{ij\theta}$ 

$$w^{n+1} - w^n + k \frac{e^{i\theta} - e^{-i\theta}}{2h} w^n - \frac{1}{Pe} k \frac{e^{-i\theta} - 2 + e^{i\theta}}{h^2} = 0$$
$$w^{n+1} = w^n \underbrace{\left(-Ci\sin\theta + \frac{2\Gamma}{Pe}\left(\cos\theta - 1\right) + 1\right)}_{A(\theta)}$$

So we may write,

$$\left|A\left(\theta\right)\right|^{2} = \left(1 - 4\frac{\Gamma}{Pe}\sin^{2}\frac{\theta}{2}\right)^{2} + C^{2}\sin^{2}\theta$$

If  $\Gamma/Pe \leq 1/2$  the first term  $\leq 1$  and since  $C^2 = k\Gamma$  then  $C^2 \sin^2 \theta \leq Pek/2 = Mk, M > 0$ , i.e.

$$\left|A\right|^2 \le 1 + Mk$$

**Theorem 2.8.1** The FC scheme is stable if  $\Gamma \leq Pe/2$  i.e. if  $|A(\theta)|^2 \leq 1 + Mk$  for some M > 0.

#### Proof

$$(Fu_h^n)^2(\xi) = A^2(h\xi) (Fu_h^{n-1})^2(\xi) = A^{2n}(h\xi) (Fu_h^0)^2(\xi) \le (1+Mk)^n (Fu_h^0)^2(\xi)$$

Therefore,

$$||Fu_{h}^{n}||_{2}^{2} \leq (1+Mk)^{n} \left| \left| Fu_{h}^{0} \right| \right|^{2} \leq (1+Mk)^{\frac{nkM}{kM}} \left| \left| Fu_{h}^{0} \right| \right|_{2}^{2} = \left( (1+Mk)^{\frac{1}{kM}} \right)^{Mt^{n}} \left| \left| Fu_{h}^{0} \right| \right|_{2}^{2} \leq e^{Mt^{n}} \left| \left| Fu_{h}^{0} \right| \right|_{2}^{2},$$

where  $t^n = nk$ .

Theorem 2.8.2 Assume:

$$\frac{\Gamma}{Pe} \le \frac{1}{2}$$

Then the solution of the FC scheme satisfies the maximum principle:

$$\max_{j} \left| u_{h,j}^{n+1} \right| \le \max_{j} \left| u_{h,j}^{n} \right|$$

for all  $n \ge 1$ , if and only if  $h \le \frac{2}{Pe}^{\dagger}$ .

<sup>&</sup>lt;sup>†</sup>Here we tacitly disregard the boundary conditions that are required since the equation involves a second derivative in space. It is quite clear, however, that the boundary conditions should be non-increasing functions of time. For example, if at the left edge of the domain we need to satisfy a Dirichlet condition that is an increasing function of time, and the approximation satisfies it exactly, then it cannot satisfy such a maximum principle in time.

**Proof** Let  $h \leq \frac{2}{Pe}$ , then:

$$\left| u_{h,j}^{n+1} \right| \le \frac{\Gamma}{Pe} \left( 1 + \frac{1}{2}hPe \right) \left| u_{h,j-1}^{n} \right| + \left( 1 - 2\Gamma Pe^{-1} \right) \left| u_{h,j}^{n} \right| + \frac{\Gamma}{Pe} \left( 1 - \frac{1}{2}hPe \right) \left| u_{h,j+1}^{n} \right| \le \max_{j} \left| u_{h,j}^{n} \right|$$

because  $h \leq \frac{2}{Pe}, \frac{\Gamma}{Pe} \leq \frac{1}{2}$ . Now, assume that

$$\max_{j} \left| u_{h,j}^{n+1} \right| \le \max_{j} \left| u_{h,j}^{n} \right|.$$

Aiming towards a contradiction, assume that  $h > \frac{2}{Pe}$ . Then, choose the initial data as follows

$$u_{h,j}^0 = \begin{cases} 1, & j = 0, 1\\ 0, & j > 1 \end{cases}$$

so that

$$u_{h,1}^{1} = \frac{\Gamma}{Pe} \left( 1 + \frac{1}{2}hPe \right) + 1 - 2\frac{\Gamma}{Pe} > \frac{\Gamma}{Pe} \left( 1 + \frac{1}{2}\frac{2}{Pe}Pe - 2 \right) + 1 = 1$$

and this gives:

$$\max_{j} \left| u_{h,j}^1 \right| > 1 = \max_{j} \left| u_{h,j}^0 \right|$$

so, by contradiction, the theorem holds.

Since violation of the maximum principle leads to unphysical solution then we must choose  $h \leq 2/Pe$  and so, if the Peclet number Pe is very large then we need to use very small h, of the order of  $Pe^{-1}$ . This corresponds to problems with boundary layers that need to be resolved by h. However, it is known from the properties of the corresponding boundary value problem that the thickness of such boundary layers is of the order of  $Pe^{-1/2}$ i.e. it can be resolved by h being of order of  $Pe^{-1/2}$ . So, this scheme requires the use of a spatial step h much less than what is needed to resolve the actual solution. The next scheme cures this problem at the expense of loss of accuracy.

#### 2.8.2 Upwinding Scheme (FB)

$$\frac{u_{h,j}^{n+1} - u_{h,j}^n}{k} + \frac{u_{h,j}^n - u_{h,j-1}^n}{h} = \frac{1}{Pe} \frac{u_{h,j-1}^n - 2u_{h,j}^n + u_{h,j+1}^n}{h^2}$$

**Theorem 2.8.3** Assume that  $\frac{\Gamma}{Pe} \leq \frac{1}{2}$ . Then the FB scheme satisfies a maximum principle provided that:

$$2\frac{\Gamma}{Pe} + C \le 1, \quad \left(C = \frac{k}{h}\right)$$

The proof is left as an exercise. This theorem implies that the maximum principle is satisfied if  $k \leq (Peh^2)/(2 + Peh)$ . Note that in the limit  $Pe \to \infty$ , this restriction on k tends to h. On the other hand, the stability condition  $k \leq Peh^2/2$  implies that if  $Pe \to \infty$  but  $Peh^2 \to const$  (i.e.  $h = O(Pe^{-1/2})$ ), k is restricted by a constant. In conclusion, if the Peclet number Pe is very large, the scheme guarantees that the solution satisfies a maximum principle, similar to the exact solution, and is stable if  $k \leq h$ , but k can be of the order of h. And this is true even if  $h = O(Pe^{-1/2})$  that is sufficient to resolve boundary layers.

The scheme has a consistency error

$$\tau_{k,h}\left(t,x\right) = \mathcal{O}\left(k+h\right)$$

 $\mathbf{as}$ 

$$\delta_t^+ u_j^n = \frac{\partial u}{\partial t} \left( t_n, x_j \right) + \mathcal{O} \left( k \right)$$

and

$$\delta_x^- u_j^n = \frac{\partial u}{\partial x} \left( t_n, x_j \right) - \frac{h}{2} \frac{\partial^2 u}{\partial x^2} \left( t_n, x_j \right) + \mathcal{O} \left( h^2 \right)$$

#### 2.8. ADVECTION-DIFFUSION EQUATION

and

$$\delta_x^2 u_j^n = \frac{\partial^2 u}{\partial x^2} \left( t_n, x_j \right) + \mathcal{O}\left( h^2 \right)$$

Combining these results we have:

$$\delta_t^+ u_j^n + \delta_x^- u_j^n - \frac{1}{Pe} \delta_x^2 u_j^n = \frac{\partial u}{\partial t} + \frac{\partial u}{\partial x} - \left(\frac{1}{Pe} + \frac{h}{2}\right) \frac{\partial^2 u}{\partial x^2} + \mathcal{O}\left(k + h^2\right)$$

i.e. up to terms that are  $O(k + h^2)$ , the scheme is consistent with the equation

$$\frac{\partial u}{\partial t} + \frac{\partial u}{\partial x} - \left(\frac{1}{Pe} + \frac{h}{2}\right)\frac{\partial^2 u}{\partial x^2} = 0.$$

This means that the scheme adds an extra diffusion term that is of order of h. If the leading order term in the spatial error is proportional to a derivative of an even order, we call the scheme dissipative since such schemes usually dissipate sharp changes (large gradients) of the solution. If the leading order term of the error is proportional to a derivative of an odd order, the scheme is called dispersive. This is because odd derivatives in PDEs lead to the so called dispersion in the solution. In terms of a Fourier decomposition of the solution, this effect is manifested in the fact that due to the presence of odd derivatives different Fourier modes travel with a different speed.

# Chapter 3

# **Introduction to Finite Elements**

# 3.1 Weighted Residual Methods

#### 3.1.1 Sobolev spaces

Sobolev space:

$$\mathbb{H}^{m}(\Omega) \equiv \left\{ v : \Omega \to \mathbb{R} : \int_{\Omega} \left( v^{2} + \left( v^{(1)} \right)^{2} + \ldots + \left( v^{(m)} \right)^{2} \right) < \infty \right\}$$

then

$$\mathbb{H}^{0}\left(\Omega\right) = L^{2}\left(\Omega\right)$$

and  $\mathbb{H}^{m}(\Omega)$  is a Hilbert space:

$$(u,v)_m = \int_{\Omega} \left( uv + u^{(1)}v^{(1)} + \ldots + u^{(m)}v^{(m)} \right) dx$$

and

$$||u||_{m} = \sqrt{\int_{\Omega} \left(u^{2} + \ldots + \left(u^{(m)}\right)^{2}\right) dx}$$

**Proposition 3.1.1** *1.*  $\mathbb{H}^{m+1}(\Omega) \subset \mathbb{H}^m(\Omega)$ 

2. If  $v \in \mathbb{H}^m(\Omega)$  then:

$$||v||_{m}^{2} = ||v||_{0}^{2} + ||v^{(1)}||_{0}^{2} + \ldots + ||v^{(m)}||_{0}^{2}$$

||v||<sub>m+1</sub> ≥ ||v||<sub>m</sub> and ||v||<sub>m+1</sub> ≥ ||v<sup>(1)</sup>||<sub>m</sub>.
 If v, w ∈ ℍ<sup>m</sup> (Ω) then |(v, w)<sub>m</sub>| ≤ ||v||<sub>m</sub> ||v||<sub>m</sub>, this is the Cauchy-Schwartz inequality.
 If v, w ∈ ℍ<sup>m</sup> (Ω) then

 $||v+w||_m \leq ||v||_m + ||w||_m$ 

**Proof** We will prove the Cauchy-Schwartz inequality

$$|(u,v)_{m}| \leq ||u||_{m} \, ||v||_{m}$$

for  $u, v \in \mathbb{H}^{m}(\Omega)$ . For any  $s \in \mathbb{R}$  we have

$$0 \le (u - sv, u - sv)_m = ||u||_m^2 + s^2 ||v||_m^2 - 2s (u, v)_m$$
43

Choose

so that:

 $s = \frac{(u, v)_m}{\left|\left|v\right|\right|_m^2}$ 

$$0 \le (u - sv, u - sv)_m = ||u||_m^2 + \frac{(u, v)_m^2}{||v||_m^2} - 2\frac{(u, v)_m^2}{||v||_m^2} = ||u||_m^2 - \frac{(u, v)_m^2}{||v||_m^2}$$

and, upon rearranging, we have:

 $\left(u,v\right)_{m}\leq\left|\left|u\right|\right|_{m}\left|\left|v\right|\right|_{m}$ 

Lu = f, in  $\Omega$ 

#### 3.1.2 Weighted Residual Formulations

Consider

with 
$$u = 0$$
, on  $\partial \Omega$ 

where

$$L = -\frac{d^2}{dx^2}$$

so that:

$$-\frac{d^2u}{dx^2} = f \quad \text{in } \Omega = (0,1)$$

u = 0, on x = 0, 1

with

Define the "trial function" space, for example it can be chosen to be

$$U = \left\{ u : u \in \mathbb{H}^2(\Omega), u = 0 \text{ on } \partial\Omega \right\}$$

and we introduce the "weight functions" (also called test) space, for this choice of U it can be chosen to be

$$W = \{ w : w \in L^2(\Omega), w = 0 \text{ on } \partial\Omega \}$$

Now, find  $u \in U$  such that

$$\int_{\Omega} \left( \frac{d^2 u}{dx^2} + f \right) w \, d\Omega = 0, \quad f \in L^2\left(\Omega\right)$$

for all  $w \in W$ . This is one "weighted residual formulation" of the original problem. This particular formulation is also usually called a strong formulation. It still defines a continuous problem and in order to discretize it we need to discretize the corresponding functional spaces i.e. to define appropriate approximations for each function in them. We select a subset  $U_h$  of U:

$$U_h \subset U, \qquad U_h = \operatorname{span} \{\phi_0, \dots, \phi_{n-1}\}$$

and  $W_h$ , a subset of W:

$$W_h \subset W, \qquad W_h = \operatorname{span} \{\psi_0, \dots, \psi_{n-1}\}$$

and then, search for the discrete approximation  $u_h$  in the form:

$$u_h = \sum_{i=0}^{n-1} c_i \phi_i$$

so that:

$$\int_{\Omega} \left( \frac{d^2}{dx^2} \sum_{i=0}^{n-1} c_i \phi_i - f \right) \psi_j \, d\Omega = 0 \quad \Rightarrow \quad \sum_{i=0}^{n-1} c_i \int_{\Omega} \frac{d^2 \phi_i}{dx^2} \psi_j = \int_{\Omega} f \psi_j \, d\Omega, \quad j = 0, 1, \dots, n-1$$

This is a linear algebraic system Lc = F where:

$$L_{ij} = \int_{\Omega} \frac{d^2 \phi_i}{dx^2} \psi_j \, d\Omega, \qquad F_i = \int_{\Omega} f \psi_i \, d\Omega$$

44

#### 3.1.3 Collocation Methods

Let

$$\psi_j\left(x\right) = \delta\left(x - x_j\right)$$

so that:  $L_{ij} = \frac{d^2 \phi_i}{dx^2} (x_j)$  and  $F_j = f(x_j)$ . Note that  $\psi_j$  are not in  $L^2$  however the formulation still makes sense if  $\phi_i \in C^0$  because the integrals are well defined. Next, we may approximate  $\frac{d^2}{dx^2} \sim \delta_x^2$ , and this gives us a finite difference method or we may select  $U_h$  as the space spanned by certain sets of orthogonal polynomials (Legendre, Chebyshev etc.), and this gives us spectral collocation.

# 3.2 Weak Methods

By far, the most popular weighted residual methods are based on the so called weak formulation of the classical problem. To obtain it we start from

$$\int_{\Omega} \left( \frac{d^2 u}{dx^2} + f \right) w \, d\Omega = 0 \quad \forall \ w \in W$$

and choose the discrete test functions space to be

$$W = \left\{ w : w \in \mathbb{H}^1_0(\Omega) \right\}.$$

Here we define

$$\mathbb{H}_{0}^{1}(\Omega) = \left\{ u \in \mathbb{H}^{1}(\Omega), \ u = 0 \text{ on } \partial \Omega \right\}.$$

Then,

$$-\int_{\Omega} \frac{du}{dx} \frac{dw}{dx} + \underbrace{\int_{\partial\Omega} \frac{du}{dx} w \, ds}_{=0} = -\int_{\Omega} f w \, d\Omega$$

Then the problem is reformulated as: Find  $u \in U$  such that:

$$\int_{\Omega} \frac{du}{dx} \frac{dw}{dx} = \int_{\Omega} f w \, d\Omega$$

A natural choice for U which makes all integrals well defined and incorporates the Dirichlet boundary conditions in the solution is U = W. Such formulation is called a *Galerkin formulation*. The space U can be discretized by means of piecewise polynomial functions, orthogonal polynomials, trigonometric polynomials, spline functions etc. These choices yield various weak methods. In the remainder of the notes we will focus on the piecewise polynomial approximations which yield the so called finite element methods.

Before we proceed with the discretization we shall prove some important results for the continuous Galerkin formulation: Find  $u \in U = \mathbb{H}^1_0(\Omega)$  such that

$$\int_{\Omega} u'v' \, d\Omega = \int_{\Omega} fv \, d\Omega, \quad \forall \, v \in U.$$

**Theorem 3.2.1** If u is a solution of the classical formulation

$$-u'' = f \quad in \ \Omega$$
$$u = 0 \quad on \ \partial\Omega$$

then u is a solution of the Galerkin formulation. If u is a solution to the Galerkin Formulation then it is a solution to the classical formulation if  $u \in C^2([0,1])$  and  $f \in C^0([0,1])$ .

**Proof** Suppose that u solves the differential formulation:

$$-u''=f$$

then take  $v \in U$  so that:

$$-\int_{\Omega} u'' v \, d\Omega = \int_{\Omega} f v \, d\Omega$$
$$\int_{\Omega} u' v' \, d\Omega = \int_{\Omega} f v \, d\Omega$$
$$(u', v')_0 = (f, v)_0$$

Suppose that u solves the Galerkin formulation so that:

$$(u',v') = (f,v)$$

then

$$\int_{\Omega} \left( u'v' - fv \right) \, d\Omega = 0 \quad \Rightarrow \quad \int_{\Omega} \left( u'' + f \right) v \, d\Omega = 0$$

Assume towards contradiction that  $u'' + f \neq 0$  in some point  $x \in \Omega$ . However,  $u'' \in \mathcal{C}^0(\Omega)$  and  $f \in \mathcal{C}^0(\Omega)$ . Therefore,

$$u'' + f \in \mathcal{C}^0(\Omega) \Longrightarrow u'' + f \neq 0$$

in an entire open interval contained in  $\Omega$ . Therefore, there exists  $(x_0, x_1) \subset \Omega$  such that u'' + f > 0 for all  $x \in (x_0, x_1)$  or u'' + f < 0 for all  $x \in (x_0, x_1)$ . We consider the first possibility and the second one can be considered in exactly the same way. Let us choose

$$v = \begin{cases} -(x-x_0)(x-x_1) & \text{in } (x_0,x_1) \\ 0 & \text{otherwise.} \end{cases}$$

For this choice it is clear that  $v \in \mathbb{H}^1_0(\Omega)$  and  $v \ge 0$  in  $\Omega$ . Now since u'' + f > 0 in  $(x_0, x_1)$  then

$$(u''+f,v)>0$$

but this contradicts the fact that u is a solution to (u'' + f, v) = 0 for all  $v \in U$ . Therefore, u'' + f = 0.

**Example** Consider the deformation of a rod from its equilibrium position under a given load f:



Then

$$E(u) = \frac{1}{2} (u', u')_0 - (f, u)_0$$

is the total energy of the system.

Therefore, E is called the energy functional of any system described by a second order elliptic differential equation.

**Theorem 3.2.2**  $u \in \mathbb{H}_0^1(\Omega)$  is a solution to the Galerkin formulation if and only if u minimizes E(v) over  $\mathbb{H}_0^1(\Omega)$ . That is,  $E(u) \leq E(v)$  for all  $v \in \mathbb{H}_0^1(\Omega)$ .

#### 3.2. WEAK METHODS

**Proof** 1. Suppose that u is a solution to the Galerkin formulation so that:  $(u', v')_0 = (f, v)_0$  for all  $v \in \mathbb{H}^1_0(\Omega)$  so that:

$$E(v) = E(u+w) = \frac{1}{2} \left( (u+w)', (u+w)' \right)_0 - (f,u+w)_0 = \frac{1}{2} \left( u', u' \right)_0 - (f,u)_0 + \frac{1}{2} \left( w', w' \right)_0 - (f,w)_0 + (u',w')_0 - (f,w)_0 + (u',w')_0 - (f,w)_0 + (u',w')_0 - (f,w)_0 + (g',w')_0 - (g',w$$

Since

we have

$$(u', w')_0 = (f, w)_0$$

$$E(v) = E(u) + \frac{1}{2} ||w'||_0^2 \ge E(u)$$

2. Suppose that  $E(u) \leq E(v)$  for all  $v \in \mathbb{H}_0^1(\Omega)$ . Then if we take  $s \in \mathbb{R}$ , E(u + sw) has a min at s = 0, and therefore:

$$\left. \frac{d}{ds} E\left( u + sw \right) \right|_{s=0} = 0$$

which we expand as:

$$\frac{d}{ds} \left( \frac{1}{2} \left( (u+sw)', (u+sw)' \right)_0 - (f,u+sw)_0 \right) \bigg|_{s=0} = 0$$
$$\frac{d}{ds} \left( E(u) + s(u',w')_0 - s(f,w)_0 + \frac{1}{2} s^2(w',w') \right) \bigg|_{s=0} = (u',w')_0 - (f,w)_0 + \underbrace{s(w',w')}_{=0} \bigg|_{s=0} = 0$$

So that:

$$(u', w')_0 = (f, w)_0$$

for all  $w \in \mathbb{H}_0^1(\Omega)$ . Therefore, u solves the Galerkin Formulation.

Consider a Hilbert space V equipped with the product  $(\cdot, \cdot)_V$  i.e. V is complete with respect to the norm  $||\cdot||_V$  induced by this product. Now, consider the mapping  $a: V \times V \to \mathbb{R}$  (further called a bilinear form) such that:

- 1.  $a(\alpha u + \beta v, w) = \alpha a(u, w) + \beta a(v, w),$
- 2.  $a(w, \alpha u + \beta v) = \alpha a(w, u) + \beta a(w, v),$
- 3. there exists  $\beta > 0$  such that  $|a(u,v)| \le \beta ||u||_V ||v||_V$  (such bilinear form is called bounded w.r.t.  $||.||_V$ ), and
- 4. there exists  $\rho > 0$  such that

 $a\left(u,u\right) \ge \rho \left|\left|u\right|\right|_{V}^{2}$ 

(such bilinear form is called coercive w.r.t.  $||.||_V$ )

Now, suppose that G(v) is a functional such that:

- 1.  $G(\alpha u + \beta v) = \alpha G(u) + \beta G(v)$  (such a functional is called linear) and
- 2. there exists  $\delta > 0$  such that  $|G(u)| \leq \delta ||u||_{V}$  (such functional is called bounded).

The following theorem is one version of a very well known result from functional analysis.

**Theorem 3.2.3 (Riesz Theorem)** If V is a Hilbert space with a given inner product  $(.,.)_*$  and if G(v) is a bounded linear functional on V (w.r.t. the norm induced by its inner product), then there is a unique element  $\hat{u} \in V$  such that  $(\hat{u}, v)_* = G(v)$  for all  $v \in V$ .

This theorem almost directly yields the proof of a somewhat restricted version of the following fundamental result in PDE theory and their numerical analysis.

**Lemma 3.2.4 (Lax-Milgram)** If a(u, v) and G(v) satisfy the six conditions stated above then there exists a unique solution  $\hat{u} \in V$  for the following problem: Find  $u \in V$  s.t.

$$a(u,v) = G(v) \quad \forall v \in V$$

If a(u, v) = a(v, u) (symmetric form) then  $\hat{u}$  is the minimizer of

$$E(v) = \frac{1}{2}a(v,v) - G(v)$$

**Proof** We shall consider only the case where a(u, v) = a(v, u).

• The proof of the first claim is a direct consequence of the Riesz theorem if we prove that a(u, v) defines an inner product on V. We know that:

$$a\left(u,u\right) \ge \rho \left|\left|u\right|\right|_{V}^{2}$$

which guarantees that

$$\sqrt{a\left(u,u\right)} = ||u||_{a}$$

is a norm on V. It is not difficult then to verify that a(u, v) defines an inner product on V. Also, it is straightforward to show that the norm  $||.||_a$  is equivalent to  $||.||_V$  i.e.:

$$\sqrt{\rho} \left| \left| u \right| \right|_V \le \left| \left| u \right| \right|_a \le \sqrt{\beta} \left| \left| u \right| \right|_V.$$

Therefore, since V is complete with respect to  $||\cdot||_V$  then it is complete with respect to the norm  $||\cdot||_a$  i.e. V is a Hilbert space with respect to the product a(.,.). As G(v) is bounded with respect to  $||\cdot||_V$ , G(v) is bounded with respect to  $||\cdot||_a$ . Since a(.,.) is an inner product on V and G(v) is bounded w.r.t. the norm induced by this inner product, the Riesz theorem implies the first claim of the lemma.

• Next we show that  $\hat{u}$  minimizes  $E(v) = \frac{1}{2}a(v,v) - G(v)$ .

$$E\left(v\right) = E\left(\hat{u} + w\right) = E\left(\hat{u}\right) + 1/2 \underbrace{a\left(w, w\right)}_{>0} \ge E\left(\hat{u}\right)$$

#### Example

$$\begin{cases} -\nabla^2 u = f & \text{in } \Omega\\ u = 0 & \text{on } \partial \Omega \end{cases}$$

Its Galerkin formulation reads: Find  $u \in \mathbb{H}^1_0(\Omega)$  s.t.:

$$(\nabla u, \nabla v)_0 = (f, v)_0 \quad \forall v \in \mathbb{H}^1_0(\Omega).$$
(3.1)

Below we prove the so-called Poincaré inequality which will guarantee that  $a(u, v) = (\nabla u, \nabla v)_0$  defines a bilinear form satisfying the conditions of the Lax-Milgram lemma.

**Lemma 3.2.5 (Poincaré Inequality)** If  $\Omega$  is a bounded domain and  $v \in \mathbb{H}_0^1(\Omega)$  then  $\exists C > 0$ , depending only on the domain  $\Omega$ , such that

$$||v||_0^2 \le C ||\nabla v||_0^2$$
.

**Proof** Consider first  $v \in \mathbb{C}^1_0(\Omega)$  i.e. v is a continuously differentiable function defined on  $\Omega$  that vanishes on its boundary. We can always find  $a \in \mathbb{R}$  large enough so that the cube  $Q = \{x \in \mathbb{R}^n : |x_j| < a, 1 \le j \le n\}$  contains  $\Omega$ . Integrating by parts in the  $x_1$  direction and taking into account that the surface integral vanishes since v = 0 on  $\partial\Omega$  we obtain:

$$\begin{aligned} ||v||_0^2 &= \int_{\Omega} v^2 dx = \int_{\Omega} 1 \cdot v^2 dx = -\int_{\Omega} x_1 \frac{\partial v^2}{\partial x_1} dx \\ &= -2 \int_{\Omega} x_1 v \frac{\partial v}{\partial x_1} dx \le 2a \int_{\Omega} |v| \left| \frac{\partial v}{\partial x_1} \right| dx. \end{aligned}$$

#### 3.3. FINITE ELEMENT METHOD (FEM)

Using the Cauchy-Schwarz inequality for the  $L^2$  inner product on  $\Omega$ , we obtain:

$$||v||_0^2 \le 2a \int_{\Omega} |v| \left| \frac{\partial v}{\partial x_1} \right| dx \le 2a ||v||_0 || \frac{\partial v}{\partial x_1} ||_0 \le 2a ||v||_0 ||\nabla v||_0.$$

Dividing by  $||v||_0$  gives the desired result with  $C = (2a)^2$  and  $v \in \mathbb{C}^1_0(\Omega)$ . Due to some classical results in functional analysis (see for example Ern, A. and Guermond, J.-L., Theory and Practice of Finite Elements, Applied Mathematical Sciences, v. 159, Springer, 2004, p. 485), we can claim that for each  $v \in \mathbb{H}^1_0(\Omega)$  we can choose a sequence of functions in  $\{v_k\}_{k=1}^{\infty} \subset \mathbb{C}^1_0(\Omega)$  such that the sequence converges to v in the  $\mathbb{H}^1\Omega$ ) norm i.e.  $||v - v_k||_0 \leq ||v - v_k||_1 \to 0$  as  $k \to \infty$  and similarly  $||\nabla v - \nabla v_k||_0 \leq ||v - v_k||_1 \to 0$ . Using the triangular inequality in the form  $||u - w||_0 \geq ||u||_0 - ||w||_0$  yields that  $||v_k||_0 \to ||v||_0$  and  $||\nabla v_k||_0 \to ||\nabla v||_0$ . This allows us to take the limit in the Poincaré inequality for  $v_k \in \mathbb{C}^1_0(\Omega)$  to derive the inequality for any function in  $\mathbb{H}^1_0(\Omega)$ .

Using this inequality we easily prove the following proposition.

**Proposition 3.2.6**  $a(u,v) = (\nabla u, \nabla v)_0$  is a bilinear, bounded, and coercive form in  $H_0^1(\Omega)$ ;  $G(v) = (f, v)_0$  is bounded in  $H_0^1(\Omega)$  if  $||f||_0 < \infty$ .

#### Proof

$$|a(u,v)| = \left| \int_{\Omega} \nabla u \nabla v dx \right| \le ||\nabla u||_0 \, ||\nabla v||_0 \le ||u||_1 \, ||v||_1$$

as  $||u||_1 = ||\nabla u||_0 + ||u||_0 \ge ||\nabla u||_0$ . Then,

$$a(u,u) = (\nabla u, \nabla u) = ||\nabla u||_0^2 = \frac{1}{2} \left( ||\nabla u||_0^2 + ||\nabla u||_0^2 \right) \ge \frac{1}{2} \left( ||u||_0^2 + ||\nabla u||_0^2 \right) = \frac{1}{2} ||u||_1^2$$

The last proposition guarantees that the bilinear form  $a(u,v) = \int_{\Omega} \nabla u \nabla v dx$  and the functional  $G(v) = (f,v)_0$  satisfy the conditions of the Lax-Milgram lemma if  $||f||_0 < \infty$ . This automatically guarantees that the corresponding weak formulation (3.1) has a unique solution in  $H_0^1(\Omega)$ .

## 3.3 Finite Element Method (FEM)

Let us define a grid:

$$\Delta = \{x_0, \ldots, x_M\}$$

and the following linear space of continuous piecewise linear functions on  $\Delta$ :

$$M_0^1(\Delta) = \operatorname{span} \{l_0, l_1, \dots, l_M\} \in \mathbb{H}^1(\Omega)$$

where

$$l_{i} = \begin{cases} \frac{x - x_{i-1}}{x_{i} - x_{i-1}}, & x_{i-1} \le x \le x_{i} \\ \frac{x - x_{i+1}}{x_{i} - x_{i+1}}, & x_{i} \le x \le x_{i+1} \\ 0 & \text{otherwise.} \end{cases}$$

Define the discrete space:

$$V_{h} = \left\{ v \in M_{0}^{1}(\Delta), \ v(0) = v(1) = 0 \right\} \subset \mathbb{H}_{0}^{1}(\Omega)$$

which is equivalent to

$$V_h = \operatorname{span} \{l_1, \ldots, l_{M-1}\}$$

Let us now consider the following discrete problem: Find  $u_h \in V_h$  such that:

$$(u'_h, v'_h)_0 = (f, v_h)_0$$

#### CHAPTER 3. INTRODUCTION TO FINITE ELEMENTS

for all  $v_h \in V_h$ . It is clear that:

$$u_h = \sum_{j=1}^{M-1} u_j l_j$$

and

$$(u'_h, l'_i)_0 = (f, l_i)_0, \quad 1 \le i \le M - 1$$

so that:

$$\left(\sum_{j=1}^{M-1} u_j l'_j, l'_i\right)_0 = (f, l_i)_0 \quad \Rightarrow \quad \sum_{j=1}^{M-1} u_j \left(l'_j, l'_i\right)_0 = (f, l_i)_0, \quad 1 \le i \le M-1.$$

The last set of equations clearly constitute a linear algebraic system:

 $Au_h = F$ 

where  $A_{ij} = (l'_i, l'_j)_0$ ,  $F_i = (f, l_i)_0$ , and  $u_{h,i} = u_i$ . Now, as

$$0 \le (v'_h, v'_h)_0 = \left(\sum_{i=1}^{M-1} v_i l'_i, \sum_{j=1}^{M-1} v_j l'_j\right)_0 = \sum_{i=1}^{M-1} \sum_{j=1}^{M-1} v_i \left(l'_i, l'_j\right)_0 v_j = v^T A v_i \left(l'_j, l'_j\right)_0 v_j =$$

the Lax-Milgram lemma guarantees that this system has a unique solution if  $f(x) \in L^2(\Omega)$ .

## 3.4 Gaussian Quadrature

We may approximate the integral of function  $\phi(x)$  over the range  $x \in (-1, 1)$  using Gaussian integration:

$$\int_{-1}^{1}\phi\left(x\right)\,dx\approx\sum_{i=1}^{n}A_{i}^{\left(n\right)}\phi\left(x_{i}^{\left(n\right)}\right)$$

where  $A_i^{(n)}$  are the weights and  $x_i^{(n)}$  are the nodes of the quadrature. n is a positive integer that controls the accuracy. The weights and nodes are determined from the condition that the quadrature is exact for all polynomials of the highest possible degree. We have 2n undetermined coefficients (n weights and n nodes) and so we can make the quadrature exact for polynomials of 2n - 1 degree (that have 2n coefficients).

The following table contains the weights and Gauss points for the first fourGauss quadratures:

n	$A_i^{(n)}$	$x_i^{(n)}$	2n - 1
1	2	0	1
2	1, 1	$-\frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}}$	3
3	$\frac{5}{9}, \frac{8}{9}, \frac{5}{9}$	$-\sqrt{\frac{3}{5}}, 0, \sqrt{\frac{3}{5}}$	5
4	0.347854, 0.652145, 0.347854, 0.652145	-0.861136, -0.339981, 0.339981, 0.861136	7

**Remark** We may approximate the integral over arbitrary bounds a and b using Gaussian Quadrature by applying the substitution:



50

#### 3.5. ERROR ESTIMATES

$$y\left(x\right) = \frac{2}{b-a}x + \frac{a+b}{a-b}$$

So that y(a) = -1, y(b) = 1, and  $dx = \frac{b-a}{2} dy$ . Therefore:

$$\int_{a}^{b} \phi(x) \, dx = \frac{b-a}{2} \int_{-1}^{1} \phi\left(\frac{(b-a)y+a+b}{2}\right) \, dy$$

# **3.5** Error Estimates

**Theorem 3.5.1** If  $u_h$  is the finite element approximation to the exact solution u of the boundary value problem

$$-u'' = f,$$
  $u(0) = u(1) = 0$ 

then there exists k > 0 such that:

$$||u - u_h||_1 \le k ||u - v_h||_1$$

for all  $v_h \in V_h$ .

**Proof** u is the exact solution, so that

for all  $w_h \in V_h$ , and

$$a\left(u_h, w_h\right) = \left(f, w_h\right)$$

 $a\left(u, w_h\right) = \left(f, w_h\right)$ 

So that

$$a\left(\underbrace{u-u_h}_{\epsilon_h}, w_h\right) = 0 \quad \Rightarrow \quad a\left(\epsilon_h, w_h\right) = 0$$

Then,

$$\frac{1}{2} ||\epsilon_h||_1^2 \le a(\epsilon_h, \epsilon_h) + a(\epsilon_h, w_h) \le a(\epsilon_h, \epsilon_h + w_h) = a(\epsilon_h, u - u_h + w_h) = a(\epsilon_h, u - v_h) \le ||\epsilon_h||_1 ||u - v_h||_1 ||u$$

So that

$$||\epsilon_h||_1 \le 2 ||u - v_h||_1$$

**Corollary 3.5.2** If u and  $u_h$  are as in the theorem, then we have

 $||u - u_h||_1 \le ch$ 

where c > 0 is independent of h.

**Proof** For the time being we will assume that  $u \in C^2(\Omega)$  i.e. u is a solution to the classical problem for  $f \in C^0(\Omega)$ . Let  $\hat{u}$  be the piecewise linear interpolant of u in  $V_h$  i.e.  $\hat{u} = \sum_{i=1}^{M-1} u_i l_i$ , where  $u_i$  are the values of u in the nodes of the grid  $x_i$ . Using Taylor expansion it is possible to prove the following estimate for the interpolation error:

$$||u' - \hat{u}'||_{L^{\infty}} \le ||u''||_{L^{\infty}} h$$

Recall that:

$$||u - u_h||_1^2 \le 4 ||u - v||_1^2$$

for all  $v \in V_h$  and as  $\hat{u} \in V_h$  we have:

$$||u - u_h||_1^2 \le 4 ||u - \hat{u}||_1^2 = 4 \left( ||u - \hat{u}||_0^2 + ||u' - \hat{u}'||_0^2 \right) \le 8 ||u' - \hat{u}'||_0^2$$

where, in the last step, we applied the Poincaré inequality. Continuing, we have:

$$||u - u_h||_1^2 \le 8 \int_0^1 (u' - \hat{u}')^2 \, dx \le 8 \, ||u' - \hat{u}'||_{L^{\infty}}^2 \le 8 \, ||u''||_{L^{\infty}}^2 \, h^2$$

where we have used the interpolation estimate  $||u' - \hat{u}'||_{L^{\infty}}^2 \leq ||u''||_{L^{\infty}}^2 h^2$ . This yields

$$||u - u_h||_1 \le \sqrt{8} ||u''||_{L^{\infty}} h = \mathcal{O}(h)$$

Applying Poincaré's inequality, we also have that:

$$||u - u_h||_1^2 = ||u - u_h||_0^2 + ||u' - u'_h||_0^2 \ge 2 ||u - u_h||_0^2$$

so that we obtain the estimate in the  $L^2$  norm:

$$||u - u_h||_0 \le ch ||u''||_{L^{\infty}}$$

This is a suboptimal estimate because from interpolation theory we know that u can be approximated with a piecewise linear function with a second order accuracy in the  $L^2$  norm.

# 3.6 Optimal Error Estimates

**Lemma 3.6.1** If  $\hat{u}$  is the finite element interpolant of  $u \in \mathbb{H}^2(\Omega)$ , i.e.  $\hat{u} = \sum_{j=1}^{M-1} u(x_j)l_j$ , then:

- 1.  $||u \hat{u}||_0 \le \left(\frac{h}{\pi}\right)^2 ||u''||_0$
- 2.  $||u' \hat{u}'||_0 \le \frac{h}{\pi} ||u''||_0$

**Proof** The domain is subdivided into elements at points:

$$x_0 = 0, x_1, x_2, \dots, x_{i-1}, x_i, \dots, x_{M-1}, x_M = 1$$

Consider the difference

 $u - \hat{u}$  in  $[0, h] = [x_0, x_1]$ 

and define  $\eta(x) = u - \hat{u} \in \mathbb{H}^1([0,h])$ . Clearly,

$$\eta\left(0\right) = \eta\left(h\right) = 0.$$

Since  $\eta \in C^0[0,h]$  it can be expanded in a uniformly convergent Fourier sine series:

$$\eta\left(x\right) = \sum_{n=1}^{\infty} \eta_n \sin\left(\frac{n\pi x}{h}\right)$$

Then, using the Parseval identity,

$$\int_0^h \eta^2 \, dx = ||\eta||_{L^2[0,h]}^2 = \frac{h}{2} \sum_{n=1}^\infty \eta_n^2.$$

Differentiating term wise we get:

$$\eta'(x) = \sum_{n=1}^{\infty} \eta_n \frac{n\pi}{h} \cos\left(\frac{n\pi x}{h}\right) \quad \Rightarrow \quad \left|\left|\eta'\right|\right|_{L^2[0,h]}^2 = \frac{h}{2} \sum_{n=1}^{\infty} \left(\eta_n \frac{n\pi}{h}\right)^2$$

and differentiating again we have:

$$\eta''(x) = -\sum_{n=1}^{\infty} \eta_n \frac{n^2 \pi^2}{h^2} \sin\left(\frac{n\pi x}{h}\right) \quad \Rightarrow \quad ||\eta''||_{L^2[0,h]}^2 = \frac{h}{2} \sum_{n=1}^{\infty} \eta_n^2 \left(\frac{n\pi}{h}\right)^4$$

So that

$$||\eta'||_{L^{2}[0,h]}^{2} = \frac{h}{2} \sum_{n=1}^{\infty} \eta_{n}^{2} \left(\frac{n\pi}{h}\right)^{2} \left(\frac{n\pi}{h}\right)^{2} \frac{h^{2}}{n^{2}\pi^{2}} \le \frac{h^{2}}{\pi^{2}} \frac{h}{2} \sum_{n=1}^{\infty} \eta_{n}^{2} \left(\frac{n\pi}{h}\right)^{4} = \frac{h^{2}}{\pi^{2}} ||\eta''||_{L^{2}[0,h]}^{2} = \frac{h^{2}}{\pi^{2}} ||u''||_{L^{2}[0,h]}^{2} = \frac{h^{2}}{\pi^{2}} |$$

and

$$||\eta||_{L^{2}[0,h]}^{2} = \frac{h}{2} \sum_{n=1}^{\infty} \eta_{n}^{2} \le \frac{h^{2}}{\pi^{2}} \frac{h}{2} \sum_{n=1}^{\infty} \eta_{n}^{2} \left(\frac{n\pi}{h}\right)^{2} = \frac{h^{2}}{\pi^{2}} \left||\eta'||_{L^{2}[0,h]}^{2} \le \left(\frac{h}{\pi}\right)^{4} \left||u''||_{L^{2}[0,h]}^{2} \le \frac{h^{2}}{\pi^{2}} \frac{h^{2}}{\pi^$$

Applying the same for any subsequent subinterval and summing the resulting inequalities we complete the proof.

**Corollary 3.6.2** 1.  $||u - \hat{u}||_1 \le \frac{h\sqrt{2}}{\pi} ||u''||_0$ 

2.  $||u - u_h||_1 \le \frac{\sqrt{8h}}{\pi} ||u''||_0$ 

**Proof** 1.

$$||u - \hat{u}||_{1}^{2} = ||u - \hat{u}||_{0}^{2} + ||u' - \hat{u}'||_{0}^{2} \le 2 ||u' - \hat{u}'||_{0}^{2} \le 2 \left(\frac{h}{\pi}\right)^{2} ||u''||_{0}^{2}$$

2. for all  $v_h \in V_h$  in the finite element space we have:

$$||u - u_h||_1 \le 2 ||u - v_h||_1$$

Take  $v_h = \hat{u}$ , then:

$$||u - u_h||_1 \le 2 ||u - \hat{u}||_1 \le \frac{h\sqrt{8}}{\pi} ||u''||_0$$

So that

$$||u - u_h||_1 \le Ch ||u''||_0$$

where C > 0 is independent of h.

**Theorem 3.6.3 (** $L^2$  **lifting theorem)** If  $u_h \in V_h$  is the Galerkin finite element approximation of:

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

then there exists  $\Gamma > 0$  such that

$$||u - u_h||_0 \le \Gamma h^2 ||u''||_0$$

**Proof** Define:

then,

$$(u', v'_h)_0 = (f, v_h)_0$$

 $\epsilon_h = u - u_h$ 

for all  $v_h \in V_h$  or,

$$a\left(u,v_{h}\right) = \left(u',v_{h}'\right)_{0}$$

so that

$$\begin{array}{l} a(u, v_h) = (f, v_h)_0 \\ a(u_h, v_h) = (f, v_h)_0 \end{array} \right\} \quad a(\epsilon_h, v_h) = 0$$

$$(3.2)$$

Now, consider the auxiliary problem

$$\begin{cases} -\phi'' = \epsilon_h, \quad 0 < x < 1\\ \phi(0) = \phi(1) = 0 \end{cases}$$

Then,

$$|\epsilon_{h}||_{0}^{2} = (\epsilon_{h}, \epsilon_{h})_{0} = -(\epsilon_{h}, \phi'')_{0} = (\epsilon'_{h}, \phi')_{0} = a(\epsilon_{h}, \phi)$$

Now, taking  $v_h = \hat{\phi}$  in (3.2), we get:

$$a\left(\epsilon_{h},\hat{\phi}\right)=0,$$

and subtracting it from the above equation we obtain:

$$\begin{aligned} ||\epsilon_{h}||_{0}^{2} &= a\left(\epsilon_{h},\phi\right) = a\left(\epsilon_{h},\phi\right) - a\left(\epsilon_{h},\hat{\phi}\right) = a\left(\epsilon_{h},\phi-\hat{\phi}\right) = \left(\epsilon_{h}',\phi'-\hat{\phi}'\right)_{0} = \int_{0}^{1} \epsilon'\left(\phi'-\hat{\phi}'\right) \, dx \\ &\leq ||\epsilon_{h}'||_{0} \left|\left|\phi'-\hat{\phi}'\right|\right|_{0} \leq ||\epsilon_{h}'||_{0} \left|\left|\phi-\hat{\phi}\right|\right|_{1} \leq ||\epsilon_{h}'||_{0} \frac{\sqrt{2}h}{\pi} \, ||\phi''||_{0} = ||\epsilon_{h}'||_{0} \frac{h\sqrt{2}}{\pi} \, ||\epsilon_{h}||_{0} \\ &\leq ||\epsilon_{h}||_{1} \frac{h\sqrt{2}}{\pi} \, ||\epsilon_{h}||_{0} \leq \Gamma h^{2} \, ||u''||_{0} \, ||\epsilon_{h}||_{0} \end{aligned}$$

So that

$$\left\|\epsilon_{h}\right\|_{0} \leq \Gamma h^{2} \left\|u''\right\|_{0}$$

for some  $\Gamma > 0$ , independent of h.

#### 3.6.1 Other Boundary Conditions

Let us consider an elliptic problem with more general boundary conditions:

$$\begin{cases} -u''(x) = f, \quad 0 < x < 1\\ u(0) = \beta_1, \quad u'(1) = \beta_2. \end{cases}$$

Since we need to satisfy non-homogeneous boundary conditions of Dirichlet and Neumann type, this time we discretize the solution with the expansion:

$$u_{h} = \beta_{1} l_{0} (x) + \sum_{j=1}^{M} u_{j} l_{j} (x)$$
(3.3)

so that

$$u_h\left(0\right) = \beta_1 l_0\left(0\right) = \beta_1.$$

Now, consider the original equation:

$$-u'' = f$$

multiply it by  $v \in V = \mathbb{H}^1(\Omega)$ , and integrate over  $\Omega$  to obtain:

$$(-u'',v)_0 = (f,v)_0 \quad \Rightarrow \quad (u',v')_0 + u'(0)v(0) - u'(1)v(1) = (f,v)_0.$$

If we approximate u with (3.3) and take the test functions to be  $v_h \in V_h = \text{span}\{l_1, \ldots, l_M\}$ , we have:

$$(u'_h, l'_j)_0 - u'_h(1) l_j(1) = (f, l_j)_0, \quad j = 1, \dots, M,$$

or

$$\left(\beta_{1}l'_{0} + \sum_{n=1}^{M} u_{n}l'_{n}, l'_{j}\right)_{0} - \beta_{2}l_{j}(1) = (f, l_{j})_{0}.$$

This yields:

$$\left(\beta_1 l'_0 + \sum_{n=1}^M u_n l'_n, l'_1\right)_0 = (f, l_1)_0,$$

$$\left(\sum_{n=1}^M u_n l'_n, l'_j\right)_0 = (f, l_j)$$

for  $j = 2, \ldots, M - 1$  and finally,

$$\left(\sum_{n=1}^{M} u_n l'_n, l'_M\right)_0 - \beta_2 = (f, l_M)_0.$$

# 3.7 Transient Problems

Consider the Initial Boundary Value Problem given by:

$$\begin{cases} \frac{\partial u}{\partial t} - \frac{\partial}{\partial x} \left[ a\left(x\right) \frac{\partial u}{\partial x} \right] = f\left(x\right) & \text{for} \quad (t, x) \in (0, T] \times (0, 1) \\ u\left(t, 0\right) = u\left(t, 1\right) = 0 & \text{for} \quad t > 0 \\ u\left(0, x\right) = g\left(x\right) & \text{for} \quad x \in (0, 1) \end{cases}$$

For  $a(x) \in \mathcal{C}^{1}(\overline{\Omega})$  such that  $0 < \alpha \leq a(x) \leq A < \infty$  and  $||f(x)||_{0} \leq L < \infty$ .

The Galerkin formulation of the problem is given by: Find  $u \in \mathbb{H}_0^1(\Omega)$  such that

$$\left(\frac{\partial u}{\partial t},v\right)_{0}+b\left(u,v\right)=\left(f\left(x\right),v\right)_{0}$$

for all  $v \in \mathbb{H}_0^1(\Omega)$ , where  $b(u, v) = \left(a(x)\frac{\partial u}{\partial x}, \frac{\partial v}{\partial x}\right)_0$ . Note that b(u, v) is a coercive and bounded bilinear form on  $H_0^1(\Omega)$  (why?). In order to discretize it we define

$$V_h = \text{span} \{ l_1, l_2, \dots l_{M-1} \}$$

and search for  $u_{h}(t, x) = \sum_{j=1}^{M-1} u_{j}(t) l_{j}(x) \in V_{h}$  such that:

$$\left(\frac{\partial u_h}{\partial t}, v_h\right)_0 + b\left(u_h, v_h\right) = (f, v_h)_0$$

for all  $v_h \in V_h$ . Then, taking into account that  $l_j$  form a basis of  $V_h$ , it is sufficient to take  $v_h = l_j, j = 1, \ldots, M-1$  and we obtain the linear system

$$\frac{\partial}{\partial t} \left( \mathbf{M} \mathbf{u}_h \right) + \mathbf{S} \mathbf{u}_h = \mathbf{F}$$

where

$$A_{ij} = \underbrace{\int_{\Omega} (l_i, l_j) \, dx}_{\text{mass matrix } \mathbf{M}} + \underbrace{\int_{\Omega} a(x) \frac{\partial l_i}{\partial x} \frac{\partial l_j}{\partial x} dx}_{\text{stiffness } \mathbf{s}}$$

for each  $i, j = 1, \ldots, M - 1$ , and

$$F_j = \int_{\Omega} fl_j \, dx, \, j = 1, \dots, M - 1; \quad \mathbf{u}_h = (u_1(t), \dots, u_{M-1}(t))^T.$$

Then, as  $\mathbf{M}$  is time-independent,

$$\mathbf{M}\frac{\partial}{\partial t}\left(\mathbf{u}_{h}\right)+\mathbf{S}\mathbf{u}_{h}=\mathbf{F}$$

We may discretize this ODE system using for example a backward difference scheme (a good choice for such problems, as we know from the previous section) to obtain the final linear system that yields the solution at time level n:

$$\mathbf{M}\frac{\mathbf{u}_h^n - \mathbf{u}_h^{n-1}}{k} + \mathbf{S}\mathbf{u}_h^n = \mathbf{F}^n.$$

Define

$$\mathbf{A} = \left(\frac{\mathbf{M}}{k} + \mathbf{S}\right),\,$$

where :

$$M_{ij} = \int_{\Omega} l_i l_j$$

and

$$S_{ij} = \int_{\Omega} a(x) \frac{\partial l_i}{\partial x} \frac{\partial l_j}{\partial x} dx \ge \alpha \int_{\Omega} \frac{\partial l_i}{\partial x} \frac{\partial l_j}{\partial x} dx.$$

M, A are spd (why?). If u is the exact solution, then

$$\frac{\partial u}{\partial t} - \frac{\partial}{\partial x} \left( a\left( x \right) \frac{\partial u}{\partial x} \right) = f$$

Now we multiply by  $v_h \in V_h$  and integrate:

$$\left(\frac{\partial u}{\partial t}, v_h\right)_0 + b\left(u, v_h\right) = (f, v_h)_0$$

Adding and subtracting  $\left(\frac{u^n - u^{n-1}}{k}, v_h\right)_0$ , we get:

$$\left(\frac{u^{n}-u^{n-1}}{k}, v_{h}\right)_{0} + b\left(u^{n}, v_{h}\right)_{0} = (f^{n}, v_{h})_{0} - \left(\frac{\partial u^{n}}{\partial t}, v_{h}\right)_{0} + \left(\frac{u^{n}-u^{n-1}}{k}, v_{h}\right)_{0}$$

which we may rewrite as:

$$\left(\frac{u^n - u^{n-1}}{k}, v_h\right)_0 + b(u^n, v_h) = (f^n, v_h)_0 + (\tau_k, v_h)_0,$$

where  $\tau_k = -\frac{\partial \mathbf{u}^n}{\partial t} + \frac{\mathbf{u}^n - \mathbf{u}^{n-1}}{k}$  is the truncation error of the backward difference approximation of the first derivative in time.

On the other hand the finite element approximation to the solution satisfies:

$$\left(\frac{u_h^n - u_h^{n-1}}{k}, v_h\right)_0 + b\left(u_h^n, v_h\right) = (f^n, v_h)_0.$$

Letting

 $u^n - u^n_h = \epsilon^n_h$ 

and subtracting the two preceding equations equations, we get:

$$\left(\frac{\epsilon_h^n - \epsilon_h^{n-1}}{k}, v_h\right) + b\left(\epsilon_h^n, v_h\right) = \left(\tau_k, v_h\right)_0$$

Now, define  $w_h$  such that

$$b(w_h, v_h) = b(u, v_h) \tag{3.4}$$

for all  $v_h \in V_h$ .  $w_h$  is called the *elliptic projection of u onto*  $V_h$ . Then we split the error  $\epsilon_h^n$  into  $\epsilon_h^n = \eta^n + \xi^n$  where  $\eta^n = u^n - w_h^n$  is the error of the elliptic projection of u. Therefore, as we have already established for the solution of (3.4)

$$||u^n - w_h^n||_0 \le \Gamma h^2 \left| \left| \frac{\partial^2 u}{\partial x^2} \right| \right|_0$$

56

But then, as  $\frac{\partial w_h}{\partial t}$  is the elliptic projection of  $\frac{\partial u}{\partial t}$  onto  $V_h$ ,

$$\left\| \left( \frac{\partial u}{\partial t} \right)^n - \left( \frac{\partial w_h}{\partial t} \right)^n \right\|_0 \le \Gamma h^2 \left\| \frac{\partial^3 u}{\partial t \partial^2 x} \right\|_0$$

Now,

$$\left(\xi^{n} - \xi^{n-1}, v_{h}\right)_{0} + kb\left(\xi^{n}, v_{h}\right) = k\left(\tau_{k}^{n}, v_{h}\right)_{0} - \left(\eta^{n} - \eta^{n-1}, v_{h}\right)_{0} - kb\left(\eta^{n}, v_{h}\right)$$

But, as  $\xi^n \in V_h$ ,

$$\left(\xi^{n} - \xi^{n-1}, \xi^{n}\right)_{0} + kb\left(\xi^{n}, \xi^{n}\right) = k\left(\tau_{k}^{n}, \xi^{n}\right)_{0} - \left(\eta^{n} - \eta^{n-1}, \xi^{n}\right)_{0} - kb\left(\eta^{n}, \xi^{n}\right)$$

and, since since  $b(\xi^n,\xi^n) \ge 0$  and  $b(\eta^n,\xi^n) = 0$ , this leads to the inequality:

$$(\xi^n - \xi^{n-1}, \xi^n)_0 \le k (\tau_k^n, \xi^n)_0 - (\eta^n - \eta^{n-1}, \xi^n)_0$$

which we may rewrite as:

$$\left\| \xi^{n} \right\|_{0}^{2} \leq \left( \xi^{n-1}, \xi^{n} \right)_{0} + k \left( \tau_{k}^{n}, \xi^{n} \right)_{0} - \left( \eta^{n} - \eta^{n-1}, \xi^{n} \right)_{0}.$$

$$(3.5)$$

Lemma 3.7.1

$$(\xi^{n-1},\xi^n)_0 \le \frac{1}{2} ||\xi^{n-1}||_0^2 + \frac{1}{2} ||\xi^n||_0^2$$

Proof

$$0 \le \left| \left| \xi^n - \xi^{n-1} \right| \right|_0^2 = \left| \left| \xi^n \right| \right|_0^2 - 2 \left( \xi^n, \xi^{n-1} \right)_0 + \left| \left| \xi^{n-1} \right| \right|_0^2$$

Upon rearranging, this yields the desired result.

#### Theorem 3.7.2 (Young Inequality)

$$(u, v)_0 \le \alpha^2 ||u||_0^2 + \frac{1}{4\alpha^2} ||v||_0^2$$

Lemma 3.7.3

$$k(\tau_{k}^{n},\xi^{n}) \leq \frac{k^{2}}{2} \int_{(n-1)k}^{nk} \left\| \left| \frac{\partial^{2}u}{\partial t^{2}} \right| \right|_{0}^{2} dt + \frac{k}{2} \left\| \xi^{n} \right\|_{0}^{2}$$

Proof

$$k(\tau_k^n,\xi^n) \le \frac{k}{2} ||\tau_k^n||_0^2 + \frac{k}{2} ||\xi^n||_0^2$$

Now using the integral form of the truncation error:  $\tau_k^n = \frac{\partial u^n}{\partial t} - \frac{u^n - u^{n-1}}{k} = \frac{1}{k} \int_{(n-1)k}^{nk} \left[t - (n-1)k\right] \frac{\partial^2 u}{\partial t^2} dt$ , and the Cauchy-Schwartz inequality for the integral in time, we obtain:

$$\begin{aligned} ||\tau_k^n||_0^2 &= \left\| \left| \frac{\partial u^n}{\partial t} - \frac{u^n - u^{n-1}}{k} \right| \right|_0^2 = \left\| \left| \frac{1}{k} \int_{(n-1)k}^{nk} \left[ t - (n-1)k \right] \frac{\partial^2 u}{\partial t^2} \, dt \right| \right|_0^2 \\ &\leq \left\| \frac{1}{k} \sqrt{\int_{(n-1)k}^{nk} \left[ t - (n-1)k \right]^2 \, dt} \sqrt{\int_{(n-1)k}^{nk} \left( \frac{\partial^2 u}{\partial t^2} \right)^2 \, dt} \right\|_0^2 \\ &\leq \left\| \frac{\sqrt{k^3}}{k} \sqrt{\int_{(n-1)k}^{nk} \left( \frac{\partial^2 u}{\partial t^2} \right)^2 \, dt} \right\|_0^2 = k \int_\Omega \int_{(n-1)k}^{nk} \left( \frac{\partial^2 u}{\partial t^2} \right)^2 \, dt \, dx \\ &= k \int_{(n-1)k}^{nk} \left\| \frac{\partial^2 u}{\partial t^2} \right\|_0^2 \, dt \end{aligned}$$

So then,

$$k\left(\tau_{k}^{k},\xi^{n}\right)_{0} \leq \frac{k}{2}\left|\left|\tau_{k}^{n}\right|\right|_{0}^{2} + \frac{k}{2}\left|\left|\xi^{n}\right|\right|_{0}^{2} \leq \frac{k^{2}}{2}\int_{(n-1)k}^{nk}\left|\left|\frac{\partial^{2}u}{\partial t^{2}}\right|\right|^{2}\,dt + \frac{k}{2}\left|\left|\xi^{n}\right|\right|_{0}^{2}\,dt + \frac{k}{2}\left|\left|\xi^{n}\right|\right$$

Lemma 3.7.4

$$\left|\left(\eta^{n}-\eta^{n-1},\xi^{n}\right)_{0}\right| \leq \frac{\Gamma^{2}h^{4}}{2} \int_{(n-1)k}^{nk} \left|\left|\frac{\partial^{3}u}{\partial t\partial x^{2}}\right|\right|_{0}^{2} dt + \frac{k}{2} \left|\left|\xi^{n}\right|\right|_{0}^{2}$$

 $\mathbf{Proof}$  Using the Cauchy-Schwartz inequality for the integral in time and subsequently, the Young's inequality we obtain

$$\begin{split} \left| \left( \eta^n - \eta^{n-1}, \xi^n \right)_0 \right| &= \left| \left( \int_{(n-1)k}^{nk} \frac{\partial \eta}{\partial t} \, dt, \xi^n \right)_0 \right| \le \left| \left( \sqrt{\int_{(n-1)k}^{nk} \left( \frac{\partial \eta}{\partial t} \right)^2 \, dt}, \sqrt{k} \xi^n \right)_0 \right| \\ &\le \frac{1}{2} \int_{\Omega} \int_{(n-1)k}^{nk} \left( \frac{\partial \eta}{\partial t} \right)^2 \, dt + \frac{k}{2} \left| |\xi^n| |_0^2 = \frac{1}{2} \int_{(n-1)k}^{nk} \int_{\Omega} \left( \frac{\partial \eta}{\partial t} \right)^2 \, dt + \frac{k}{2} \left| |\xi^n| |_0^2 \\ &= \frac{1}{2} \int_{(n-1)k}^{nk} \left| \left| \frac{\partial \eta}{\partial t} \right| \right|_0^2 \, dt + \frac{k}{2} \left| |\xi^n| |_0^2 \end{split}$$

Then, applying the result

$$\left\| \left| \frac{\partial w_h}{\partial t} - \frac{\partial u}{\partial t} \right| \right\|_0^2 \le \Gamma^2 h^4 \left\| \left| \frac{\partial^2}{\partial x^2} \left( \frac{\partial u}{\partial t} \right) \right\|_0^2 = \Gamma^2 h^4 \left\| \left| \frac{\partial^3 u}{\partial x^2 \partial t} \right| \right\|_0^2,$$

we obtain the final claim of the lemma.

**Lemma 3.7.5** Let  $k \leq \frac{1}{4}$ , then:

$$||\xi^{n}||_{0}^{2} \leq 2\left||\xi^{0}||_{0}^{2} + 2k^{2}\left|\left|\frac{\partial^{2}u}{\partial t^{2}}\right|\right|_{L^{2}(0,T,L^{2}(\Omega))}^{2} + 2h^{4}\Gamma^{2}\left|\left|\frac{\partial^{3}u}{\partial x^{2}\partial t}\right|\right|_{L^{2}(0,T,L^{2}(\Omega))}^{2} + 4k\sum_{m=0}^{n-1}||\xi^{m}||_{0}^{2} + 2k^{2}\left||\frac{\partial^{2}u}{\partial t^{2}}\right|_{L^{2}(0,T,L^{2}(\Omega))}^{2} + 2h^{4}\Gamma^{2}\left|\frac{\partial^{3}u}{\partial x^{2}\partial t}\right|_{L^{2}(0,T,L^{2}(\Omega))}^{2} + 4k\sum_{m=0}^{n-1}||\xi^{m}||_{0}^{2} + 2k^{2}\left||\frac{\partial^{2}u}{\partial t^{2}}\right|_{L^{2}(0,T,L^{2}(\Omega))}^{2} + 2h^{4}\Gamma^{2}\left|\frac{\partial^{3}u}{\partial x^{2}\partial t}\right|_{L^{2}(0,T,L^{2}(\Omega))}^{2} + 4k\sum_{m=0}^{n-1}||\xi^{m}||_{0}^{2} + 2k^{2}\left||\frac{\partial^{2}u}{\partial t^{2}}\right|_{L^{2}(0,T,L^{2}(\Omega))}^{2} + 2h^{4}\Gamma^{2}\left|\frac{\partial^{3}u}{\partial x^{2}\partial t}\right|_{L^{2}(0,T,L^{2}(\Omega))}^{2} + 2h^{4}\Gamma^{2}\left|\frac{\partial^{3}u}{\partial x^{2}}\right|_{L^{2}(0,T,L^{2}(\Omega))}^{2} + 2h^{4}\Gamma^{2}\left|\frac{\partial^{3}u}{\partial x^{2}}\right|_{L^{2}(0,T,L^{2}(\Omega)}^{2} + 2h^{4}\Gamma^{2}\left|\frac{\partial^{3}u}{\partial x^{2}}\right|_{L^{2}(0,T,L^{2}(\Omega)}^{2} + 2h^$$

Where we define

$$||f(t,x)||_{L^{2}(0,T,L^{2}(\Omega))} = \int_{0}^{T} \int_{\Omega} f^{2}, \qquad f \in L^{2}((0,T) \times \Omega).$$

**Proof** Subbing the results of the last three lemmas into (3.5) we easily obtain:

$$\frac{1}{2} \left\| \xi^m \right\|_0^2 \le \frac{1}{2} \left\| \xi^{m-1} \right\|_0^2 + \frac{k^2}{2} \int_{(m-1)k}^{mk} \left\| \left| \frac{\partial^2 u}{\partial t^2} \right\|_0^2 dt + \frac{\Gamma^2 h^4}{2} \int_{(m-1)k}^{mk} \left\| \frac{\partial^3 u}{\partial x^2 \partial t} \right\|_0^2 dt + k \left\| \xi^m \right\|_0^2,$$

or:

$$||\xi^{m}||_{0}^{2} \leq \left|\left|\xi^{m-1}\right|\right|_{0}^{2} + k^{2} \int_{(m-1)k}^{mk} \left|\left|\frac{\partial^{2}u}{\partial t^{2}}\right|\right|_{0}^{2} dt + \Gamma^{2}h^{4} \int_{(m-1)k}^{mk} \left|\left|\frac{\partial^{3}u}{\partial x^{2}\partial t}\right|\right|_{0}^{2} dt + 2k \left||\xi^{m}\right||_{0}^{2}.$$

Now we sum for  $m = 1, \ldots, n$  to get:

$$\begin{split} ||\xi^{n}||_{0}^{2} &\leq \left|\left|\xi^{0}\right|\right|_{0}^{2} + k^{2} \int_{0}^{nk} \left\|\frac{\partial^{2}u}{\partial t^{2}}\right\|_{0}^{2} dt + \Gamma^{2}h^{4} \int_{0}^{nk} \left\|\frac{\partial^{3}u}{\partial t\partial x^{2}}\right\|_{0}^{2} dt + 2k \sum_{m=1}^{n-1} ||\xi^{m}||_{0}^{2} + 2k \, ||\xi^{n}||_{0}^{2} dt \\ (1-2k) \, ||\xi^{n}||_{0}^{2} &\leq \left|\left|\xi^{0}\right|\right|_{0}^{2} + k^{2} \int_{0}^{nk} \left\|\frac{\partial^{2}u}{\partial t^{2}}\right\|_{0}^{2} + \Gamma^{2}h^{4} \int_{0}^{nk} \left\|\frac{\partial^{3}u}{\partial x^{2}\partial t}\right\|_{0}^{2} dt + 2k \sum_{m=1}^{n-1} ||\xi^{m}||_{0}^{2} . \end{split}$$

#### 3.7. TRANSIENT PROBLEMS

Since  $k \leq 1/4$  then  $1 - 2k \geq 1/2$  we obtain:

$$||\xi^{n}||_{0}^{2} \leq 2\left|\left|\xi^{0}\right|\right|_{0}^{2} + 2k^{2} \int_{0}^{nk} \left|\left|\frac{\partial^{2}u}{\partial t^{2}}\right|\right|_{0}^{2} + 2\Gamma^{2}h^{4} \int_{0}^{nk} \left|\left|\frac{\partial^{3}u}{\partial x^{2}\partial t}\right|\right|_{0}^{2} dt + 4k \sum_{m=0}^{n-1} \left|\left|\xi^{m}\right|\right|_{0}^{2} dt + 2k^{2} \int_{0}^{n} \left|\left|\frac{\partial^{2}u}{\partial t^{2}}\right|\right|_{0}^{2} dt + 2k^{2} \int_{0}^{n} \left|\frac{\partial^{2}u}{\partial t^{2}}\right|^{2} dt + 2k^{2} \int_{0}^{$$

and this competes the proof.

**Lemma 3.7.6 (Gronwall Lemma)** If the sequence  $\{\xi^n\}_{n=1}^N$  satisfies

$$|\xi^n| \le \nu + \mu \sum_{m=0}^{n-1} |\xi^m|, \qquad \nu, \mu \ge 0$$

then,

$$|\xi^{n}| \le e^{\mu N} \left(\nu + \mu \left|\xi^{0}\right|\right), \quad \forall n \in \{1, \dots, N\}$$

**Proof** Consider the sequence  $z_n$  defined by

$$z^{n} = \nu + \mu \sum_{m=0}^{n} |\xi^{m}| \quad \Rightarrow \quad z^{n} - z^{n-1} = \mu |\xi^{n}|$$

so that,

$$z^{n} \leq (1+\mu) z^{n-1}$$
  
$$z^{n-1} \leq (1+\mu) z^{n-2} \leq \ldots \leq (1+\mu)^{\frac{(n-1)\mu}{\mu}} z^{0} \leq e^{\mu(n-1)} z^{0} \leq e^{\mu N} z^{0}$$

because:

$$(1+\mu)^{\frac{1}{\mu}} \le e$$

Then,

$$|\xi^{n}| \le z^{n-1} \le e^{\mu N} z^{0} = e^{\mu N} \left(\nu + \mu \left|\xi^{0}\right|\right)$$

for n = 1, ..., N.

Now let  $\mu = 4k$  and

$$\nu = 2 \left| \left| \xi^0 \right| \right|_0^2 + 2k^2 \left| \left| \frac{\partial^2 u}{\partial t^2} \right| \right|_{L^2(0,T;L^2(\Omega))}^2 + 2\Gamma^2 h^4 \left| \left| \frac{\partial^3 u}{\partial t \partial x^2} \right| \right|_{L^2(0,T;L^2(\Omega))}^2$$

Denoting the final time instant by T = nk and using the results of lemma 3.7.5 we find that

$$||\xi^{n}||_{0}^{2} \leq e^{4T} \left( 3 \left| \left| \xi^{0} \right| \right|_{0}^{2} + 2k^{2} \left| \left| \frac{\partial^{2}u}{\partial t^{2}} \right| \right|_{L^{2}(0,T;L^{2}(\Omega))}^{2} + 2h^{4}\Gamma^{2} \left| \left| \frac{\partial^{3}u}{\partial t\partial x^{2}} \right| \right|_{L^{2}(0,T;L^{2}(\Omega))}^{2} \right)$$

Let's choose  $u_h^0 = w_h^0$  i..e.

$$b\left(u_{h}^{0}, v_{h}\right) = b\left(g, v_{h}\right)$$

for all  $v_h \in V_h$ . Then  $\xi^0 = 0$ . Therefore,

$$\left|\left|\xi^{n}\right|\right|_{0}^{2} \leq e^{4T} \left(2k^{2} \left|\left|\frac{\partial^{2}u}{\partial t^{2}}\right|\right|_{L^{2}(0,T;L^{2}(\Omega))}^{2} + 2\Gamma^{2}h^{4} \left|\left|\frac{\partial^{3}u}{\partial t\partial x^{2}}\right|\right|_{L^{2}(0,T;L^{2}(\Omega))}^{2}\right) = \mathcal{O}\left(k^{2} + h^{4}\right)$$

**Theorem 3.7.7** Consider the scheme:

$$\mathbf{M}\frac{\mathbf{u}^n-\mathbf{u}^{n-1}}{k}+\mathbf{A}\mathbf{u}^n=\mathbf{f}$$

If  $k < \frac{1}{4}$  and  $u_h^0 = w_h^0$ , then

$$\left|\left|\epsilon_{h}^{n}\right|\right|_{0} \leq \mathcal{O}\left(k+h^{2}\right)$$

Proof

$$||\xi^n||_0^2 \le a^2k^2 + b^2h^4 \le \left(ak + bh^2\right)^2$$

where we define:

$$\begin{aligned} a^{2} &= 2e^{4T} \left| \left| \frac{\partial^{2}u}{\partial t^{2}} \right| \right|_{L^{2}(0,T;L^{2}(\Omega))}^{2} \\ b^{2} &= 2e^{4T}\Gamma^{2} \left| \left| \frac{\partial^{3}u}{\partial t\partial x^{2}} \right| \right|_{L^{2}(0,T;L^{2}(\Omega))}^{2} \end{aligned}$$

Then,

$$||\epsilon_h^n||_0 \leq ||\eta^n||_0 + ||\xi^n||_0 \leq \sqrt{2}e^{2T} \left| \left| \frac{\partial^2 u}{\partial t^2} \right| \right|_{L^2\left(0,T;L^2(\Omega)\right)} k + h^2 \Gamma \left[ \max_{0 \leq n \leq N} \left| \left| \frac{\partial^2 u^n}{\partial x^2} \right| \right|_0 + \sqrt{2}e^{2T} \left| \left| \frac{\partial^3 u}{\partial t \partial x^2} \right| \right|_{L^2\left(0,T;L^2(\Omega)\right)} \right] \leq \frac{1}{2} \left| \frac{\partial^2 u^n}{\partial t^2} \right| \left| \frac{\partial^2 u$$

# 3.8 Finite Elements in Two Dimensions

Consider the model problem:



This is the global numbering, the local numbering (for i = 1, 2, 3) is:



where k is the element number and  $i_j^k$  is the global number of the j-th local point in the k-th element.

#### 3.8.1 Finite Element Basis Functions

With each global node i we associate one basis function  $\phi_i(x, y)$ . If  $N_j$  is the j-th point in the grid, we require:

$$\phi_i\left(N_j\right) = \delta_{ij}$$

where  $\left.\phi_{j}\right|_{e_{k}} \in P^{1}$ , the space of all first degree polynomials. Then,

$$V_h = \operatorname{span} \{\phi_1, \dots, \phi_M\}$$

where M is the total number of points in the grid.



where  $i_r^k = N_i$  and r can be 1, 2, or 3 (local numbers).



Now, we define:

$$T_k \equiv \begin{cases} x = x \, (\xi, \eta) = a_1 \xi + a_2 \eta + a_3 \\ y = y \, (\xi, \eta) = b_1 \xi + b_2 \eta + b_3 \end{cases}$$

Then,

$$\phi_{(\underline{1})}^{k}(\xi,\eta) = 1 - \xi - \eta, \quad \phi_{(\underline{2})}^{k}(\xi,\eta) = \xi, \quad \phi_{(\underline{3})}^{k}(\xi,\eta) = \eta$$

So that:

$$\int_{\Omega} \phi_i(x, y) \phi_j(x, y) \ d\Omega = \sum_{k=1}^{K_e} \int_{e_k} \phi_i(\xi, \eta) \phi_j(\xi, \eta) \ de_k$$

Where  $K_e$  is the total number of finite elements,  $de_k = |\mathcal{J}_{T_k}| d\xi d\eta$ , and  $\mathcal{J}_{T_k}$  is the Jacobian of  $T_k$ .

#### 3.8.2 Discrete Galerkin Formulation

Find  $u_h \in V_h$  such that:

$$\int_{\Omega} p u_{h,x} v_{h,x} + \int_{\Omega} p u_{h,y} v_{h,y} + \int_{\Omega} q u_h v_h = \int_{\Omega} g v_h$$

and, for simplicity, we choose  $\gamma \equiv 0$ , p = constant and q = constant. Then, substitutiting

$$u_h = \sum_{i=1}^M u_i \phi_i$$

and  $v_h = \phi_j$  for j = 1, ..., M yields the linear system:

$$(p\mathbf{S}+q\mathbf{M})\mathbf{u}_h=\mathbf{G}$$

where:

$$S_{ij} = p \int_{\Omega} \nabla \phi_i \cdot \nabla \phi_j \, d\Omega \quad \text{(the stiffness matrix)}$$
$$M_{ij} = q \int_{\Omega} \phi_i \phi_j \, d\Omega \quad \text{(the mass matrix)}$$