

Scientific Computing WS 2019/2020

Lecture 21

Jürgen Fuhrmann

[juergen.fuhrmann@wias-berlin.de](mailto:juergen.fuhrmann@wias-berlin.de)

## Convection-Diffusion problems

# The convection - diffusion equation

Search function  $u : \Omega \rightarrow \mathbb{R}$  such that

$$\begin{aligned} -\nabla \cdot (D\vec{\nabla}u - u\vec{v}) &= f \quad \text{in } \Omega \\ u &= g \quad \text{on } \Gamma \end{aligned}$$

- $u(x)$ : species concentration, temperature
- $\vec{j} = D\vec{\nabla}u - u\vec{v}$ : species flux
- $D$ : diffusion coefficient
- $\vec{v}(x)$ : velocity of medium (e.g. fluid)
  - Given analytically
  - Solution of free flow problem (Navier-Stokes equation)
  - Flow in porous medium (Darcy equation):  $\vec{v} = -\kappa\vec{\nabla}p$  where

$$-\nabla \cdot (\kappa\vec{\nabla}p) = 0$$

- For constant density, the divergence condition  $\nabla \cdot \vec{v} = 0$  holds.

## Weak formulation

- Let  $u_g \in H^1(\Omega)$  a lifting of  $g$ . Find  $u \in H^1(\Omega)$  such that

$$u = u_g + \phi$$

$$\int_{\Omega} (D\vec{\nabla}\phi - \phi\vec{v}) \cdot \vec{\nabla}w \, d\vec{x} = \int_{\Omega} fw \, d\vec{x} + \int_{\Omega} (D\vec{\nabla}u_g - u_g\vec{v}) \cdot \vec{\nabla}w \quad \forall w \in H_0^1(\Omega)$$

Is this bilinear form coercive? - Use Lax-Milgram, for it being true, is not necessary that the bilinear form is self-adjoint.

- It follows, that

$$\int_{\Omega} (D\vec{\nabla}u - u\vec{v}) \cdot \vec{\nabla}w \, d\vec{x} = \int_{\Omega} fw \, d\vec{x} \quad \forall w \in H_0^1(\Omega)$$

- Green's theorem: If  $w = 0$  on  $\partial\Omega$ :

$$\int_{\Omega} \vec{v} \cdot \vec{\nabla}w \, d\vec{x} = - \int_{\Omega} w \nabla \cdot \vec{v} \, d\vec{x}$$

## Coercivity of bilinear form

Regard the convection contribution to the coercivity estimate:

$$-\int_{\Omega} u\vec{v} \cdot \vec{\nabla} u \, d\vec{x} = \int_{\Omega} u\vec{\nabla} \cdot (u\vec{v}) \, d\vec{x} \quad \text{Green's theorem}$$

$$\int_{\Omega} u^2 \vec{\nabla} \cdot \vec{v} \, d\vec{x} + \int_{\Omega} u\vec{v} \cdot \vec{\nabla} u \, d\vec{x} = \int_{\Omega} u\vec{\nabla} \cdot (u\vec{v}) \, d\vec{x} \quad \text{Product rule}$$

$$\int_{\Omega} u^2 \vec{\nabla} \cdot \vec{v} \, d\vec{x} + 2 \int_{\Omega} u\vec{v} \cdot \vec{\nabla} u \, d\vec{x} = 0 \quad \text{Equation difference}$$

$$\int_{\Omega} u\vec{v} \cdot \vec{\nabla} u \, d\vec{x} = 0 \quad \text{Divergence condition } \vec{\nabla} \cdot \vec{v} = 0$$

Then

$$\int_{\Omega} (D\vec{\nabla} u - u\vec{v}) \cdot \vec{\nabla} u \, d\vec{x} = \int_{\Omega} D\vec{\nabla} u \cdot \vec{\nabla} u \, d\vec{x} \geq C \|u\|_{H_0^1(\Omega)}$$

One could allow for fixed sign of  $\nabla \cdot \vec{v}$ .

## The Lax-Milgram lemma

**Theorem:** Let  $V$  be a Hilbert space. Let  $a : V \times V \rightarrow \mathbb{R}$  be a self-adjoint bilinear form, and  $f$  a linear functional on  $V$ . Assume  $a$  is coercive, i.e.

$$\exists \alpha > 0 : \forall u \in V, a(u, u) \geq \alpha \|u\|_V^2.$$

Then the problem: find  $u \in V$  such that

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

admits one and only one solution with an a priori estimate

$$\|u\|_V \leq \frac{1}{\alpha} \|f\|_{V'}$$



## Convection diffusion problem: maximum principle

- Let  $f \leq 0$ ,  $\nabla \cdot \vec{v} = 0$
- Let  $g^\sharp = \sup_{\partial\Omega} g$ .
- Let  $w = (u - g^\sharp)^+ = \max\{u - g^\sharp, 0\} \in H_0^1(\Omega)$
- Consequently,  $w \geq 0$
- As  $\vec{\nabla} u = \vec{\nabla}(u - g^\sharp)$  and  $\vec{\nabla} w = 0$  where  $w \neq u - g^\sharp$ , one has

$$0 \geq \int_{\Omega} fw \, d\vec{x} = \int_{\Omega} D(\vec{\nabla} u - u\vec{v})\vec{\nabla} w \, d\vec{x} \quad \text{Variational identity}$$

$$= \int_{\Omega} D(\vec{\nabla} w - w\vec{v})\vec{\nabla} w \, d\vec{x} - Dg^\sharp \int_{\Omega} \vec{v} \cdot \vec{\nabla} w \, d\vec{x} \quad \text{Replace } u \text{ by } w$$

$$= \int_{\Omega} D(\vec{\nabla} w - w\vec{v})\vec{\nabla} w \, d\vec{x} + Dg^\sharp \int_{\Omega} w\vec{\nabla} \cdot \vec{v} \, d\vec{x} \quad \text{Green}$$

$$\geq C\|w\|_{H_0^1(\Omega)} \quad \text{Coercivity, } \nabla \cdot \vec{v} = 0$$

- Therefore:  $w = (u - g^\sharp)^+ = 0$  and  $u \leq g^\sharp$
- Similar for minimum part

## Mimimax for convection-diffusion

**Theorem:** If  $\nabla \cdot \vec{v} = 0$ , the weak solution of the inhomogeneous Dirichlet problem

$$\begin{aligned} -\nabla \cdot (D\vec{\nabla}u - u\vec{v}) &= f & \text{in } \Omega \\ u &= g & \text{on } \partial\Omega \end{aligned}$$

fulfills the global minimax principle: it attains its maximum at the boundary if  $f \leq 0$  and attains its minimum at the boundary if  $f \geq 0$ .

**Corollary:** If  $f = 0$  then  $u$  attains both its minimum and its maximum at the boundary.

**Corollary:** Local minimax principle: This is true of any subdomain  $\omega \subset \Omega$ .



## Finite volumes for convection diffusion

$$\begin{aligned} -\nabla \cdot \vec{j} &= 0 \quad \text{in } \Omega \\ \vec{j} \cdot \vec{n} + \alpha u &= g \quad \text{on } \Gamma = \partial\Omega \end{aligned}$$

- Integrate time discrete equation over control volume

$$\begin{aligned} 0 &= - \int_{\omega_k} \nabla \cdot \vec{j} d\omega = - \int_{\partial\omega_k} \vec{j} \cdot \vec{n}_k d\gamma \\ &= - \sum_{I \in \mathcal{N}_k} \int_{\sigma_{kl}} \vec{j} \cdot \vec{n}_{kl} d\gamma - \int_{\gamma_k} \vec{j} \cdot \vec{n} d\gamma \\ &\approx \sum_{I \in \mathcal{N}_k} \underbrace{\frac{|\sigma_{kl}|}{h_{kl}} g_{kl}(u_k, u_I)}_{\rightarrow A_\Omega} + \underbrace{|\gamma_k| \alpha u_k - |\gamma_k| g_k}_{\rightarrow A_\Gamma} \end{aligned}$$

- $A = A_\Omega + A_\Gamma$

## Central Difference Flux Approximation

- $g_{kl}$  approximates normal convective-diffusive flux between control volumes  $\omega_k, \omega_l$ :  $g_{kl}(u_k - u_l) \approx -(D\vec{\nabla}u - u\vec{v}) \cdot \vec{n}_{kl}$
- Let  $\sigma_{kl} = \omega_k \cap \omega_l$   
Let  $v_{kl} = \frac{1}{|\sigma_{kl}|} \int_{\sigma_{kl}} \vec{v} \cdot \vec{n}_{kl} d\gamma$  approximate the normal velocity  $\vec{v} \cdot \vec{n}_{kl}$
- Central difference flux:

$$\begin{aligned}g_{kl}(u_k, u_l) &= D(u_k - u_l) + h_{kl} \frac{1}{2}(u_k + u_l)v_{kl} \\ &= (D + \frac{1}{2}h_{kl}v_{kl})u_k - (D - \frac{1}{2}h_{kl}v_{kl})u_l\end{aligned}$$

- if  $v_{kl}$  is large compared to  $h_{kl}$ , the corresponding matrix (off-diagonal) entry may become positive
- Non-positive off-diagonal entries only guaranteed for  $h \rightarrow 0$  !
- If all off-diagonal entries are non-positive, we can prove the discrete maximum principle

## Simple upwind flux discretization

- Force correct sign of convective flux approximation by replacing central difference flux approximation  $h_{kl}\frac{1}{2}(u_k + u_l)v_{kl}$  by

$$\left( \begin{cases} h_{kl}u_k v_{kl}, & v_{kl} < 0 \\ h_{kl}u_l v_{kl}, & v_{kl} > 0 \end{cases} \right) = h_{kl}\frac{1}{2}(u_k + u_l)v_{kl} + \underbrace{\frac{1}{2}h_{kl}|v_{kl}|}_{\text{Artificial Diffusion } \tilde{D}} (u_k - u_l)$$

- Upwind flux:

$$\begin{aligned} g_{kl}(u_k, u_l) &= D(u_k - u_l) + \begin{cases} h_{kl}u_k v_{kl}, & v_{kl} > 0 \\ h_{kl}u_l v_{kl}, & v_{kl} < 0 \end{cases} \\ &= (D + \tilde{D})(u_k - u_l) + h_{kl}\frac{1}{2}(u_k + u_l)v_{kl} \end{aligned}$$

- M-Property guaranteed unconditionally !
- Artificial diffusion introduces error: second order approximation replaced by first order approximation

## Exponential fitting flux I

- Project equation onto edge  $x_K x_L$  of length  $h = h_{kl}$ , let  $v = -v_{kl}$ , integrate once

$$u' - uv = j$$

$$u|_0 = u_k$$

$$u|_h = u_l$$

- Linear ODE
- Solution of the homogeneous problem:

$$u' - uv = 0$$

$$u'/u = v$$

$$\ln u = u_0 + vx$$

$$u = K \exp(vx)$$

## Exponential fitting II

- Solution of the inhomogeneous problem: set  $K = K(x)$ :

$$K' \exp(vx) + vK \exp(vx) - vK \exp(vx) = -j$$

$$K' = -j \exp(-vx)$$

$$K = K_0 + \frac{1}{v}j \exp(-vx)$$

- Therefore,

$$u = K_0 \exp(vx) + \frac{1}{v}j$$

$$u_k = K_0 + \frac{1}{v}j$$

$$u_l = K_0 \exp(vh) + \frac{1}{v}j$$

## Exponential fitting III

- Use boundary conditions

$$\begin{aligned}K_0 &= \frac{u_k - u_l}{1 - \exp(vh)} \\u_k &= \frac{u_k - u_l}{1 - \exp(vh)} + \frac{1}{v}j \\j &= \frac{v}{\exp(vh) - 1}(u_k - u_l) + vu_k \\&= v \left( \frac{1}{\exp(vh) - 1} + 1 \right) u_k - \frac{v}{\exp(vh) - 1} u_l \\&= v \left( \frac{\exp(vh)}{\exp(vh) - 1} \right) u_k - \frac{v}{\exp(vh) - 1} u_l \\&= \frac{-v}{\exp(-vh) - 1} u_k - \frac{v}{\exp(vh) - 1} u_l \\&= \frac{B(-vh)u_k - B(vh)u_l}{h}\end{aligned}$$

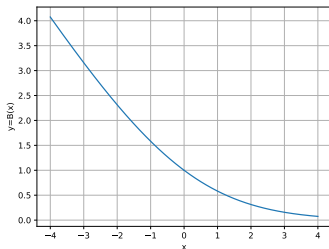
where  $B(\xi) = \frac{\xi}{\exp(\xi) - 1}$ : Bernoulli function

## Exponential fitting IV

- General case:  $Du' - uv = D(u' - u\frac{v}{D})$
- Upwind flux:

$$g_{kl}(u_k, u_l) = D(B(\frac{-v_{kl}h_{kl}}{D})u_k - B(\frac{v_{kl}h_{kl}}{D})u_l)$$

- Allen+Southwell 1955
- Scharfetter+Gummel 1969
- Ilin 1969
- Chang+Cooper 1970
- Guaranteed sign pattern,  $M$  property!



## Exponential fitting: Artificial diffusion

- Difference of exponential fitting scheme and central scheme
- Use:  $B(-x) = B(x) + x \Rightarrow$

$$B(x) + \frac{1}{2}x = B(-x) - \frac{1}{2}x = B(|x|) + \frac{1}{2}|x|$$

$$\begin{aligned}D_{art}(u_k - u_l) &= D(B(\frac{-vh}{D})u_k - B(\frac{vh}{D})u_l) - D(u_k - u_l) + h\frac{1}{2}(u_k + u_l)v \\ &= D(\frac{-vh}{2D} + B(\frac{-vh}{D}))u_k - D(\frac{vh}{2D} + B(\frac{vh}{D})u_l) - D(u_k - u_l) \\ &= D\left(\frac{1}{2}\left|\frac{vh}{D}\right| + B\left(\left|\frac{vh}{D}\right|\right) - 1\right)(u_k - u_l)\end{aligned}$$

- Further, for  $x > 0$ :

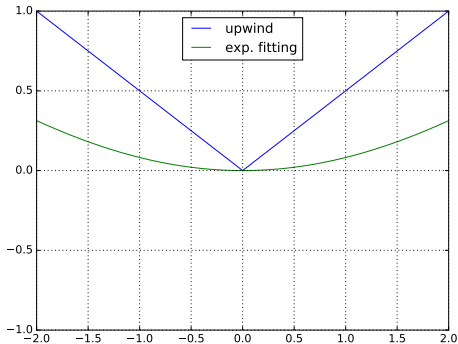
$$\frac{1}{2}x \geq \frac{1}{2}x + B(x) - 1 \geq 0$$

- Therefore

$$\frac{|vh|}{2} \geq D_{art} \geq 0$$



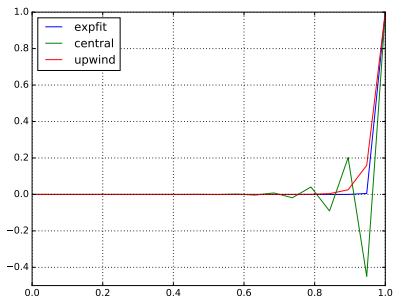
## Exponential fitting: Artificial diffusion II



Comparison of artificial diffusion functions  $\frac{1}{2}|x|$  (upwind)  
and  $\frac{1}{2}|x| + B(|x|) - 1$  (exp. fitting)

## Convection-Diffusion test problem, $N=20$

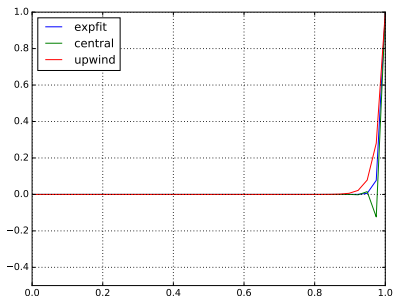
- $\Omega = (0, 1)$ ,  $-\nabla \cdot (D \vec{\nabla} u + uv) = 0$ ,  $u(0) = 0$ ,  $u(1) = 1$
- $V = 1$ ,  $D = 0.01$



- Exponential fitting: sharp boundary layer, for this problem it is exact
- Central differences: unphysical
- Upwind: larger boundary layer

## Convection-Diffusion test problem, $N=40$

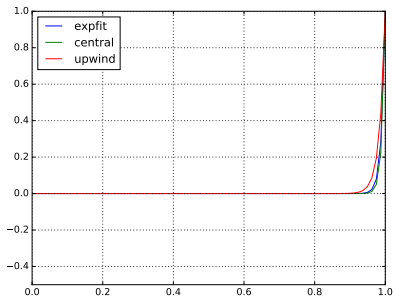
- $\Omega = (0, 1)$ ,  $-\nabla \cdot (D \vec{\nabla} u + uv) = 0$ ,  $u(0) = 0$ ,  $u(1) = 1$
- $V = 1$ ,  $D = 0.01$



- Exponential fitting: sharp boundary layer, for this problem it is exact
- Central differences: unphysical, but less “wiggles”
- Upwind: larger boundary layer

## Convection-Diffusion test problem, $N=80$

- $\Omega = (0, 1)$ ,  $-\nabla \cdot (D\vec{\nabla}u + uv) = 0$ ,  $u(0) = 0$ ,  $u(1) = 1$
- $V = 1$ ,  $D = 0.01$



- Exponential fitting: sharp boundary layer, for this problem it is exact
- Central differences: grid is fine enough to yield M-Matrix property, good approximation of boundary layer due to higher convergence order
- Upwind: “smearing” of boundary layer

## 1D convection diffusion summary

- Upwinding and exponential fitting unconditionally yield the  $M$ -property of the discretization matrix
- Exponential fitting for this case (zero right hand side, 1D) yields exact solution. It is anyway “less diffusive” as artificial diffusion is optimized
- Central scheme has higher convergence order than upwind (and exponential fitting) but on coarse grid it may lead to unphysical oscillations
- For 2/3D problems, sufficiently fine grids to stabilize central scheme may be prohibitively expensive
- Local grid refinement may help to offset artificial diffusion

# Discrete minimax principle

- $Au = f$
- $A$ : matrix from diffusion or convection- diffusion
- $A$  irreducibly diagonally dominant, positive main diagonal entries, negative off diagonal entries

$$a_{ii}u_i = \sum_{j \neq i} -a_{ij}u_j + f_i$$

$$u_i = \sum_{j \neq i, a_{ij} \neq 0} -\frac{a_{ij}}{a_{ii}}u_j + \frac{f_i}{a_{ii}}$$

- For interior points,,  $a_{ii} = -\sum_{j \neq i} a_{ij}$
- Assume  $i$  is interior point. Assume  $f_i \geq 0 \Rightarrow$

$$u_i \geq \min_{j \neq i, a_{ij} \neq 0} u_j \sum_{j \neq i, a_{ij} \neq 0} -\frac{a_{ij}}{a_{ii}} = \min_{j \neq i, a_{ij} \neq 0} u_j$$

- Assume  $i$  is interior point. Assume  $f_i \leq 0 \Rightarrow$

$$u_i \leq \max_{j \neq i, a_{ij} \neq 0} u_j \sum_{j \neq i, a_{ij} \neq 0} -\frac{a_{ij}}{a_{ii}} = \max_{j \neq i, a_{ij} \neq 0} u_j$$

## Discussion of discrete minimax principle I

- P1 finite elements, Voronoi finite volumes: matrix graph  $\equiv$  triangulation of domain
- The set  $\{j \neq i, a_{ij} \neq 0\}$  is exactly the set of neighbor nodes
- Solution in point  $x_i$  estimated by solution in neighborhood
- The estimate can be propagated to the boundary of the domain

## Discussion of discrete minimax principle II

- Minimax principle + positivity/nonnegativity of solutions can be seen as an important qualitative property of the physical process
- Along with good approximation quality, its preservation in the discretization process may be necessary
- Guaranteed for irreducibly diagonally dominant matrices
- Nonnegativity for nonnegative right hand sides guaranteed by *M*-Property
- Finite volume method may be preferred as it can guarantee these properties for boundary conforming Delaunay grids.



## Convection-diffusion and finite elements

Search function  $u : \Omega \rightarrow \mathbb{R}$  such that

$$\begin{aligned} -\vec{\nabla} \cdot (D \vec{\nabla} u - u \vec{v}) &= f \quad \text{in } \Omega \\ u &= g \quad \text{on } \partial\Omega \end{aligned}$$

- Assume  $v$  is divergence-free, i.e.  $\nabla \cdot v = 0$ .
- Then the main part of the equation can be reformulated as

$$-\vec{\nabla} \cdot (D \vec{\nabla} u) + \vec{v} \cdot \vec{\nabla} u = 0 \quad \text{in } \Omega$$

yielding a weak formulation: find  $u \in H^1(\Omega)$  such that  $u - g \in H_0^1(\Omega)$  and  $\forall w \in H_0^1(\Omega)$ ,

$$\int_{\Omega} D \vec{\nabla} u \cdot \vec{\nabla} w \, dx + \int_{\Omega} \vec{v} \cdot \vec{\nabla} u \, w \, dx = \int_{\Omega} f w \, dx$$

- Galerkin formulation: find  $u_h \in V_h$  with bc. such that  $\forall w_h \in V_h$

$$\int_{\Omega} D \vec{\nabla} u_h \cdot \vec{\nabla} w_h \, dx + \int_{\Omega} \vec{v} \cdot \vec{\nabla} u_h \, w_h \, dx = \int_{\Omega} f w_h \, dx$$

## Convection-diffusion and finite elements II

- Galerkin ansatz has similar problems as central difference ansatz in the finite volume/finite difference case  $\Rightarrow$  stabilization ?
- Most popular: streamline upwind Petrov-Galerkin (SUPG)

$$\int_{\Omega} D \vec{\nabla} u_h \cdot \vec{\nabla} w_h \, dx + \int_{\Omega} \vec{v} \cdot \vec{\nabla} u_h w_h \, dx + S(u_h, w_h) = \int_{\Omega} f w_h \, dx$$

with

$$S(u_h, w_h) = \sum_K \int_K (-\vec{\nabla} \cdot (D \vec{\nabla} u_h - u_h \vec{v}) - f) \delta_K \vec{v} \cdot w_h \, dx$$

where  $\delta_K = \frac{h_K^v}{2|\vec{v}|} \xi\left(\frac{|\vec{v}| h_K^v}{D}\right)$  with  $\xi(\alpha) = \coth(\alpha) - \frac{1}{\alpha}$  and  $h_K^v$  is the size of element  $K$  in the direction of  $\vec{v}$ .

## Convection-diffusion and finite elements III

- Many methods to stabilize, *none* guarantees M-Property even on weakly acute meshes ! (V. John, P. Knobloch, Computer Methods in Applied Mechanics and Engineering, 2007)

- Comparison paper:

M. Augustin, A. Caiazzo, A. Fiebach, J. Fuhrmann, V. John, A. Linke, and R. Umla, "An assessment of discretizations for convection-dominated convection-diffusion equations," *Comp. Meth. Appl. Mech. Engrg.*, vol. 200, pp. 3395–3409, 2011:

- if it is necessary to compute solutions without spurious oscillations: use FVM, taking care on the construction of an appropriate grid might be essential for reducing the smearing of the layers,
- if sharpness and position of layers are important and spurious oscillations can be tolerated: often the SUPG method is a good choice.

- Topic of ongoing research

## Transient problems

## Time dependent Robin boundary value problem

- Choose final time  $T > 0$ . Regard functions  $(x, t) \rightarrow \mathbb{R}$ .

$$\partial_t u - \nabla \cdot D\vec{\nabla} u = f \quad \text{in } \Omega \times [0, T]$$

$$D\vec{\nabla} u \cdot \vec{n} + \alpha u = g \quad \text{on } \partial\Omega \times [0, T]$$

$$u(x, 0) = u_0(x) \quad \text{in } \Omega$$

- This is an initial boundary value problem
- This problem has a weak formulation in the Sobolev space  $L^2([0, T], H^1(\Omega))$ , which then allows for a Galerkin approximation in a corresponding subspace
- We will proceed in a simpler manner: first, perform a finite difference discretization in time, then perform a finite element (finite volume) discretization in space.
  - *Rothe method*: first discretize in time, then in space
  - *Method of lines*: first discretize in space, get a huge ODE system, then apply perform discretization

## Time discretization

- Choose time discretization points  $0 = t^0 < t^1 \dots < t^N = T$
- let  $\tau^n = t^n - t^{n-1}$   
For  $i = 1 \dots N$ , solve

$$\frac{u^n - u^{n-1}}{\tau^n} - \nabla \cdot D\vec{\nabla} u^\theta = f \quad \text{in } \Omega \times [0, T]$$
$$D\vec{\nabla} u_\theta \cdot \vec{n} + \alpha u^\theta = g \quad \text{on } \partial\Omega \times [0, T]$$

where  $u^\theta = \theta u^n + (1 - \theta)u^{n-1}$

- $\theta = 1$ : backward (implicit) Euler method  
Solve PDE problem in each timestep. First order accuracy in time.
- $\theta = \frac{1}{2}$ : Crank-Nicolson scheme  
Solve PDE problem in each timestep. Second order accuracy in time.
- $\theta = 0$ : forward (explicit) Euler method  
First order accurate in time. This does not involve the solution of a PDE problem  $\Rightarrow$  Cheap? What do we have to pay for this ?

# Finite volumes for time dependent hom. Neumann problem

Search function  $u : \Omega \times [0, T] \rightarrow \mathbb{R}$  such that  $u(x, 0) = u_0(x)$  and

$$\partial_t u - \nabla \cdot D\vec{\nabla} u = 0 \quad \text{in } \Omega \times [0, T]$$

$$D\vec{\nabla} u \cdot \vec{n} = 0 \quad \text{on } \Gamma \times [0, T]$$

- Given control volume  $\omega_k$ , integrate equation over space-time control volume  $\omega_k \times (t^{n-1}, t^n)$ , divide by  $\tau^n$ :

$$\begin{aligned} 0 &= \int_{\omega_k} \left( \frac{1}{\tau^n} (u^n - u^{n-1}) - \nabla \cdot D\vec{\nabla} u^\theta \right) d\omega \\ &= \frac{1}{\tau} \int_{\omega_k} (u^n - u^{n-1}) d\omega - \int_{\partial\omega_k} D\vec{\nabla} u^\theta \cdot \vec{n}_k d\gamma \\ &= - \sum_{I \in \mathcal{N}_k} \int_{\sigma_{kl}} D\vec{\nabla} u^\theta \cdot \vec{n}_{kl} d\gamma - \int_{\gamma_k} D\vec{\nabla} u^\theta \cdot \vec{n} d\gamma - \frac{1}{\tau} \int_{\omega_k} (u^n - u^{n-1}) d\omega \\ &\approx \underbrace{\frac{|\omega_k|}{\tau^n} (u_k^n - u_k^{n-1})}_{\rightarrow M} + \underbrace{\sum_{I \in \mathcal{N}_k} \frac{|\sigma_{kl}|}{h_{kl}} (u_k^\theta - u_l^\theta)}_{\rightarrow A} \end{aligned}$$

# Matrix equation

- Resulting matrix equation:

$$\frac{1}{\tau^n} (Mu^n - Mu^{n-1}) + Au^\theta = 0$$

$$\frac{1}{\tau^n} Mu^n + \theta Au^n = \frac{1}{\tau^n} Mu^{n-1} + (\theta - 1)Au^{n-1}$$

$$u^n + \tau^n M^{-1} \theta Au^n = u^{n-1} + \tau^n M^{-1} (\theta - 1) Au^{n-1}$$

- $M = (m_{kl})$ ,  $A = (a_{kl})$  with

$$a_{kl} = \begin{cases} \sum_{l' \in \mathcal{N}_k} D \frac{|\sigma_{kl'}|}{h_{kl'}} & l = k \\ -D \frac{\sigma_{kl}}{h_{kl}}, & l \in \mathcal{N}_k \\ 0, & \text{else} \end{cases}$$

$$m_{kl} = \begin{cases} |\omega_k| & l = k \\ 0, & \text{else} \end{cases}$$

- $\Rightarrow \theta A + M$  is strictly diagonally dominant!



## A matrix norm estimate

**Lemma:** Assume  $A$  has positive main diagonal entries, nonpositive off-diagonal entries and row sum zero. Then,  $\|(I + A)^{-1}\|_{\infty} \leq 1$

**Proof:** Assume that  $\|(I + A)^{-1}\|_{\infty} > 1$ .  $I + A$  is a reducible  $M$ -matrix, thus  $(I + A)^{-1}$  has positive entries. Then for  $\alpha_{ij}$  being the entries of  $(I + A)^{-1}$ ,

$$\max_{i=1}^n \sum_{j=1}^n \alpha_{ij} > 1.$$

Let  $k$  be a row where the maximum is reached. Let  $e = (1 \dots 1)^T$ . Then for  $v = (I + A)^{-1}e$  we have that  $v > 0$ ,  $v_k > 1$  and  $v_k \geq v_j$  for all  $j \neq k$ . The  $k$ th equation of  $e = (I + A)v$  then looks like

$$\begin{aligned} 1 &= v_k + v_k \sum_{j \neq k} |a_{kj}| - \sum_{j \neq k} |a_{kj}| v_j \\ &\geq v_k + v_k \sum_{j \neq k} |a_{kj}| - \sum_{j \neq k} |a_{kj}| v_k \\ &= v_k \\ &> 1 \end{aligned}$$

## Stability estimate

- Matrix equation again:

$$u^n + \tau^n M^{-1} \theta A u^n = u^{n-1} + \tau^n M^{-1} (\theta - 1) A u^{n-1} =: B^n u^{n-1}$$
$$u^n = (I + \tau^n M^{-1} \theta A)^{-1} B^n u^{n-1}$$

- From the lemma we have  $\|(I + \tau^n M^{-1} \theta A)^n\|_\infty \leq 1$   
 $\Rightarrow \|u^n\|_\infty \leq \|B^n u^{n-1}\|_\infty$ .
- For the entries  $b_{kl}^n$  of  $B^n$ , we have

$$b_{kl}^n = \begin{cases} 1 + \frac{\tau^n}{m_{kk}} (\theta - 1) a_{kk}, & k = l \\ \frac{\tau^n}{m_{kk}} (\theta - 1) a_{kl}, & \text{else} \end{cases}$$

- In any case,  $b_{kl} \geq 0$  for  $k \neq l$ .  
If  $b_{kk} \geq 0$ , one estimates  $\|B\|_\infty = \max_{k=1}^N \sum_{l=1}^N b_{kl}$ .
- But

$$\sum_{l=1}^N b_{kl} = 1 + (\theta - 1) \frac{\tau^n}{m_{kk}} \left( a_{kk} + \sum_{l \in \mathcal{N}_k} a_{kl} \right) = 1 \Rightarrow \|B\|_\infty = 1.$$

## Stability conditions

- For a shape regular triangulation in  $\mathbb{R}^d$ , we can assume that  $m_{kk} = |\omega_k| \sim h^d$ , and  $a_{kl} = \frac{|\sigma_{kl}|}{h_{kl}} \sim \frac{h^{d-1}}{h} = h^{d-2}$ , thus  $\frac{a_{kk}}{m_{kk}} \leq \frac{1}{Ch^2}$
- $b_{kk} \geq 0$  gives

$$(1 - \theta) \frac{\tau^n}{m_{kk}} a_{kk} \leq 1$$

- A sufficient condition is that for some  $C > 0$ ,

$$(1 - \theta) \frac{\tau^n}{Ch^2} \leq 1$$

$$(1 - \theta)\tau^n \leq Ch^2$$

- Method stability:
  - Implicit Euler:  $\theta = 1 \Rightarrow$  unconditional stability !
  - Explicit Euler:  $\theta = 0 \Rightarrow$  CFL condition  $\tau \leq Ch^2$
  - Crank-Nicolson:  $\theta = \frac{1}{2} \Rightarrow$  CFL condition  $\tau \leq 2Ch^2$   
Tradeoff stability vs. accuracy.

## Stability discussion

- $\tau \leq Ch^2$  CFL == “Courant-Friedrichs-Levy”
- Explicit (forward) Euler method can be applied on very fast systems (GPU), with small time step comes a high accuracy in time.
- Implicit Euler: unconditional stability – helpful when stability is of utmost importance, and accuracy in time is less important
- For hyperbolic systems (pure convection without diffusion), the CFL conditions is  $\tau \leq Ch$ , thus in this case explicit computations are ubiquitous
- Comparison for a fixed size of the time interval. Assume for implicit Euler, time accuracy is less important, and the number of time steps is independent of the size of the space discretization.

	1D	2D	3D
# unknowns	$N = O(h^{-1})$	$N = O(h^{-2})$	$N = O(h^{-3})$
# steps	$M = O(N^2)$	$M = O(N)$	$M = O(N^{2/3})$
complexity	$M = O(N^3)$	$M = O(N^2)$	$M = O(N^{5/3})$

## Backward Euler: discrete maximum principle

$$\begin{aligned}\frac{1}{\tau^n} M u^n + A u^n &= \frac{1}{\tau} M u^{n-1} \\ \frac{1}{\tau^n} m_{kk} u_k^n + a_{kk} u_k^n &= \frac{1}{\tau^n} m_{kk} u_k^{n-1} + \sum_{k \neq l} (-a_{kl}) u_l^n \\ u_k^n &= \frac{1}{\frac{1}{\tau^n} m_{kk} + \sum_{l \neq k} (-a_{kl})} \left( \frac{1}{\tau^n} m_{kk} u_k^{n-1} + \sum_{l \neq k} (-a_{kl}) u_l^n \right) \\ &\leq \frac{\frac{1}{\tau^n} m_{kk} + \sum_{l \neq k} (-a_{kl})}{\frac{1}{\tau^n} m_{kk} + \sum_{l \neq k} (-a_{kl})} \max(\{u_k^{n-1}\} \cup \{u_l^n\}_{l \in \mathcal{N}_k}) \\ &\leq \max(\{u_k^{n-1}\} \cup \{u_l^n\}_{l \in \mathcal{N}_k})\end{aligned}$$

- Provided, the right hand side is zero, the solution in a given node is bounded by the value from the old timestep, and by the solution in the neighboring points.
- No new local maxima can appear during time evolution
- There is a continuous counterpart which can be derived from weak solution
- Sign pattern is crucial for the proof.

## Backward Euler: Nonnegativity

$$u^n + \tau^n M^{-1} A u^n = u^{n-1}$$

$$u^n = (I + \tau^n M^{-1} A)^{-1} u^{n-1}$$

- $(I + \tau^n M^{-1} A)$  is an M-Matrix
- If  $u_0 > 0$ , then  $u^n > 0 \forall n > 0$

# Mass conservation

- Equivalent of  $\int_{\Omega} \nabla \cdot D\vec{\nabla} u d\vec{x} = \int_{\partial\Omega} D\vec{\nabla} u \cdot \vec{n} d\gamma = 0$ :

$$\begin{aligned} \sum_{k=1}^N \left( a_{kk} u_k + \sum_{l \in \mathcal{N}_k} a_{kl} u_l \right) &= \sum_{k=1}^N \sum_{l=1, l \neq k}^N a_{kl} (u_l - u_k) \\ &= \sum_{k=1}^N \sum_{l=1, l < k}^N (a_{kl} (u_l - u_k) + a_{lk} (u_k - u_l)) \\ &= 0 \end{aligned}$$

- $\Rightarrow$  Equivalent of  $\int_{\Omega} u^n d\vec{x} = \int_{\Omega} u^{n-1} d\vec{x}$ :

- $\sum_{k=1}^N m_{kk} u_k^n = \sum_{k=1}^N m_{kk} u_k^{n-1}$

# Weak formulation of time step problem

- Weak formulation: search  $u \in H^1(\Omega)$  such that  $\forall v \in H^1(\Omega)$

$$\frac{1}{\tau^n} \int_{\Omega} u^n v \, dx + \theta \int_{\Omega} D \vec{\nabla} u^n \vec{\nabla} v \, dx =$$
$$\frac{1}{\tau^n} \int_{\Omega} u^{n-1} v \, dx + (1 - \theta) \int_{\Omega} D \vec{\nabla} u^{n-1} \vec{\nabla} v \, dx$$

- Matrix formulation

$$\frac{1}{\tau^n} M u^n + \theta A u^n = \frac{1}{\tau^n} M u^{n-1} + (1 - \theta) A u^{n-1}$$

- $M$ : mass matrix,  $A$ : stiffness matrix.
- With FEM, Mass matrix lumping important for getting the previous estimates