

Condition number of the stiffness matrix:

Assume the triangulation J_h satisfies

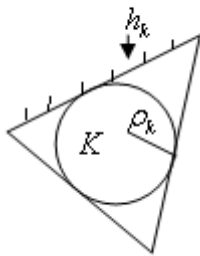
$$h_k \geq \beta_1 h \quad h = \max_{k \in J_h} h_k, \quad h_k = \text{diameter of element } k$$

$$\frac{\rho_k}{h_k} \geq \beta_2 \quad \rho_k = \text{radius of the circle inscribed in } k = \int \nabla u \nabla v dx$$

Consider bilinear form $a(u, v)$ satisfying the coercivity and continuity

on $H_{2,0}^1(\Omega)$ ($V_h \subset H_{2,0}^1$)

$$\begin{cases} (i) a(u, v) \geq \delta \|u\|^2 \\ (ii) a(u, v) \leq \beta \|u\| \|v\| \quad (\Omega \subset \mathbb{R}^2) \end{cases}$$



The following Lemma holds

Lemma 1. \exists constants c and C (depends on α, β) such that for all

$$v_h = \sum_{i=1}^m v_i \varphi_i \in V_h, \text{ the following inequalities hold.}$$

$$(10) \quad ch^2 |\vec{v}|^2 \leq \|v_h\|^2 \leq Ch^2 |\vec{v}|^2 \quad (\vec{v} = (v_1, v_2, \dots, v_m))$$

$$(11) \quad a(v_h, v_h) \equiv \int_{\Omega} |\nabla v_h|^2 dx \leq ch^{-2} \|v_h\|^2 \quad \left(\begin{array}{l} \text{inverse estimate compare with Poincaré} \\ \text{inequality } \|v\|_{L^2} < \|\nabla v\|_{L^2} \end{array} \right)$$

(prove: Skip. See also Johnson's section 7.7)

(exercise)

With the help of the Lemma 1, we can show the condition number of the stiff matrix K is $O(h^{-2})$.

$$\text{cond}(K) = \|K\| \|K^{-1}\|_{(\text{matrix norm})} \quad \|K\| = \sup_{x \in \mathbb{R}^n} \left(\frac{\|Ax\|}{\|x\|} \right)$$

Since

$$\frac{\vec{v}^T K \vec{v}}{|\vec{v}|^2} = \frac{a(v_h, v_h)}{|\vec{v}|^2} \stackrel{(11)}{\leq} ch^{-2} \frac{\|v_h\|^2}{|\vec{v}|^2} \stackrel{(10)}{\leq} c^* \sup_v \lambda_{\max} \leq c^{**} \left(\begin{array}{l} \vec{v}: \text{evector corresponding} \\ \text{to max evalue} \end{array} \right)$$

and

$$\frac{\vec{v}^T K \vec{v}}{|\vec{v}|^2} = \frac{a(v_h, v_h)}{|\vec{v}|^2} \geq \partial \frac{\|v_h\|^2}{|\vec{v}|^2} \geq c^{**} h^2 \stackrel{\text{inf}}{\Rightarrow} \lambda_{\min} \geq c^{**} h^2 \left(\begin{array}{l} \vec{v}: \text{evector corresponding} \\ \text{to min evalue} \end{array} \right)$$

$$\Rightarrow \lambda_{\max}(K) \leq C^* \text{ and } \lambda_{\min}(K) \geq c^{**} h^2$$

$$\Rightarrow \text{cond}(K) \leq \widetilde{Ch^{-2}} \left(\widetilde{C = \frac{c^*}{c^{**}}} \right)$$

(exercise: In 3D, $ch^3 |\vec{v}| \leq \|u\|^2 < ch^3 |v| \Rightarrow \text{cond}(K) < \widetilde{Ch^2}$)

Remark:

(1) Recall that when solving linear system $Ax = b$ by iterative method,

(let x be the iterative solution i.e. $Ax = b + \Delta b$), we have

$$\frac{|\Delta x|}{|x|} \leq \text{cond}(A) \frac{|r|}{|b|} \quad \text{here } r = b - Ax \text{ (the residual)}$$

If the relative error is required to be less than ε , the relative residual $\left(\frac{|r|}{|b|} \right)$

should be required to be less than $\frac{\varepsilon}{\text{cond}(A)} \stackrel{\approx}{\approx} \varepsilon \cdot h^2$
(\uparrow FEM stiff matrix)

$$\Rightarrow \frac{|\Delta x|}{|x|} < \varepsilon \Rightarrow |\Delta x| < \varepsilon |x| \stackrel{(10)}{=} \varepsilon \|x_h\| h^{-1}$$

\uparrow vector

(2) How large ε should be?

$$\text{Recall that } \|x_{\text{true}} - x_{\text{FEM}}\|_{L^2} \stackrel{\uparrow \text{Duality argument}}{<} ch^2 \|x_{\text{true}}\|_{H^2(\Omega)} \stackrel{\downarrow \text{Assuming regularity estimate}}{<} ch^2 \|f\|$$

(x_{FEM} is the FEM solution, x_{true} is the PDE solution)

Since we don't want $|\Delta x| > \|x_{\text{true}} - x_{\text{FEM}}\|$, $\left(\begin{array}{l} \text{i.e. we want } |\Delta x| \ll h^2 \|f\| \\ |\Delta x| < \varepsilon |x_{\text{FEM}}| \end{array} \right)$

($\Delta x = x_{\text{iter}} - x_{\text{FEM}}$)

$$\text{we set } \varepsilon \ll \underbrace{\frac{\|f\|}{|x_{\text{FEM}}|}}_v \cdot h^2$$

$$\Rightarrow \varepsilon \ll \underbrace{r}_{\substack{r \text{ dep on } f \\ \text{and } x_{\text{true}}}} \cdot h^3, \frac{\|f\|}{\|x_{\text{FEM}}\|} h^3 > \frac{\|f\|}{\|x_{\text{true}}\|} \cdot \frac{\|x_{\text{true}}\|}{\|x_{\text{true}}\| + ch^2 \|x_{\text{true}}\|_{H^2}} h^3$$

The a posteriori error estimation:

$$-\Delta u + cu = f$$

consider $a(u, v) = \int \nabla u \nabla v + c \int uv$ and u be the solution of

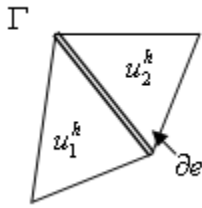
$$a(u, v) = \int f v dx = L(f) \quad \|u\|_E: \text{the energy norm}$$

$$\text{let } \begin{cases} e = u - u^h \\ e^h = u^h - u^h, \text{ here } u^h = Iu \text{ (interpolant of } u) \\ \eta = u^h - u \end{cases}$$

clearly we have $e = e^h + \eta$

Since

$$\begin{aligned}
\|e\|_E^2 &= |a(e, e)| = |a(e^h + \eta, e)| \stackrel{\substack{\text{orthogonality} \\ (a(x_h, e) = 0 \ \forall_h \in V_h)}}{=} |a(\eta, e)| \\
&= |a(\eta, u - u^h)| \\
&= |a(\eta, u) - a(\eta, u^h)| \\
&= \left| L(\eta) - \left(\int_{\Omega} \nabla \eta \nabla u^h + c \int_{\Omega} \eta u^h \right) \right| \\
&= \left| L(\eta) - \sum_{K \in J_h} \int_e \nabla \eta \nabla u^h + c \int_{\Omega} \eta u^h \right| \\
&\stackrel{\substack{\text{(integration} \\ \text{by part)}}{=}}{\left| L(\eta) + \int_{\Omega} \eta (\Delta u^h - c u^h) + \sum_{\substack{s \in \partial \\ K \in J_h}} \int_K \eta \underbrace{\left[\nabla u^h \cdot \vec{n}_1 \right]}_{(*)} \right|} \\
&\leq \underbrace{\left| \langle \eta, r^h \rangle \right|}_{(I)} + \sum_{s=1}^{\text{\# of edges}} \underbrace{\left| \int_s \eta(x) \underbrace{\left[\nabla u_k^h \cdot \vec{n}_s \right]}_{J_{k,s}^h} \right|}_{(II)} \\
r^h &= f + \Delta u_h - c u_h \\
\left((*) \text{ flux jump: } \int_{\hat{c}e_1} (\nabla u_2^h - \nabla u_1^h) \eta = \left[\underbrace{(\nabla u_2^h - \nabla u_1^h) \cdot n}_J \right] \right)
\end{aligned}$$



Using more advance interpolation estimation, it can be shown

$$\begin{cases} \|u - Iu\|_{L^2(K)} \leq c_1 h_k \|e\|_{E,K} \\ \|u - Iu\|_{L^2(\partial K)} \leq c_2 h_k^{\frac{1}{2}} \|e\|_{E,K} \end{cases}$$

\uparrow energy norm, \uparrow element norm

$$\Rightarrow (I) \leq \|h^{-1}\eta\| \|hr^h\|_{L^2} \leq c \|hr^h\|_{L^2} \|e\|_E$$

$$(II) \leq \sum_{K \in J_h} \sum_{s \rightarrow \partial K} h_k^{-\frac{1}{2}} \eta(x_k) \cdot h^{\frac{1}{2}} J_{k,s}^h$$

$$\leq \sum_{K \in J_h} \sum_{s \rightarrow \partial K} c_2 \|e\|_{E,K} \left[\int_{\partial s} \left(h^{\frac{1}{2}} J_{K,s}^h \right)^2 ds \right]^{\frac{1}{2}} \leq \sum_{k \in J_h} c_2 \overline{\|e\|_{E,K}} (hJ_{k,s}^h)$$

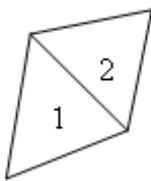
$$\leq c_2 \left(\sum_{k \in J_h} \|e\|_{E,K}^2 \right)^{\frac{1}{2}} \left(\sum_{k \in J_h} (hJ_{k,s}^h)^2 \right)^{\frac{1}{2}}$$

let $R_h = \sum_{k \in J_h} (J_{k,s}^h)$, we have

$$\|e\|_E^2 \leq \tilde{c} \left(\|hr^h\|_{L^2(\Omega)} + \sum_{\substack{K \in J_h \\ s \in \text{edges of } J_h}} h_s \cdot J_{K,s}^h \right)$$

\uparrow
 only depends on u_h

here $J_{K,s}^h = \begin{cases} \frac{1}{2} |\nabla u_1^h \cdot \vec{n}_s - \nabla u_2^h \cdot \vec{n}_s|, & s \text{ in the interior of } \Omega \\ 0 & \text{if } s \text{ on the } \partial\Omega \end{cases}$

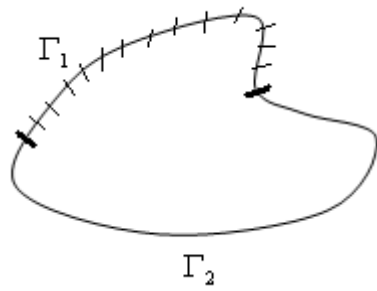


FEM for parabolic problem

consider the heat equation

$$(+)\begin{cases} \dot{u} - \text{div}(\mu \nabla u) = f & \text{in } \Omega \times I \\ u = 0 & \text{on } \Gamma_1 \times I \\ \mu \frac{\partial u}{\partial n} = 0 & \text{on } \Gamma_2 \times I \\ u(x, 0) = \dot{u}(x) & x \in \Omega \end{cases}$$

(initial boundary value problem)



1-D Model problem:

$$(+) \begin{cases} \frac{d}{dt}u - \frac{d^2u}{dx^2} = f \\ u(0,t) = u(\pi,t) = 0 \\ u(x,0) = u^0 \end{cases}$$

In case $f = 0$, by separation of variables

$$\begin{aligned}
 \text{consider } u(x,t) = e^{-w^2t}h(x) &\Rightarrow \frac{\partial u}{\partial t} = -w^2e^{-w^2t}h(x), \quad \frac{\partial^2 u}{\partial x^2} = e^{-w^2t}h''(x) \\
 &\Rightarrow -w^2e^{-w^2t}h(x) - e^{-w^2t}h''(x) = 0 \\
 &\Rightarrow e^{-w^2t}(h''(x) + w^2h(x)) = 0
 \end{aligned}$$



Consider $h(x) = de^{iwx}$ for any $d \in \mathbb{C}$

From $u(0,t) = u(\pi,t) = 0$, we have $h(x) = (-i)e^{iwx}$

real part
 $\Rightarrow h(x) = \sin wx$

Therefore a general solution of (+) has the following form

$$u(x,t) = \sum_w w e^{-w^2 t} \sin wx \quad - (**)$$

Since $u(x,0) = \sum_w c_w \sin wx = u^0 \Rightarrow c_w = \frac{2}{\pi} \int (\sin wx) u^0 dx$
↑
fourier coefficient

$$\Rightarrow u(x,t) = \sum_{w=1}^{\infty} c_w e^{-w^2 t} \sin wx$$

\Rightarrow each component $\sin wx$ lives on a time scale $o(w^{-2})$

\Rightarrow High freq modes quickly get damped! - (1)

\Rightarrow The solution u becomes smoother as $t \rightarrow \infty$

But $u(x,t)$ will not be smooth for $t \ll 1$

$$\left\| \dot{u}(t) \right\| = \left(\int_{\Omega} u^2(x,t) dx \right)^{\frac{1}{2}} = \left\| \frac{\partial^2 u}{\partial x^2} \right\| = \left\| \sum_w c_w w^2 e^{-iwt} \sin wx \right\| \rightarrow \infty \text{ as } t \rightarrow 0 \quad - (2)$$

$$\left\| \dot{u}(t) \right\| = \left\| \frac{\partial^2 u}{\partial x^2} \right\| \rightarrow 0 \text{ the rate is depends on how small } c_w \text{ is for large } w$$

↑
multitude of
high freq modes

\Rightarrow In general smooth initial u^0 gives small c_w for large w_2

An initial phase for t small where certain derivatives of u are large is called an "initial transient"

Observation: based on (1) & (2)

(i) initial transient \Rightarrow small time step for discretization of $\frac{d}{df} u$

if rough initial u^0 (oscillating discontinuity etc.)

\Rightarrow small mech size for discretization of u_{xx}

(ii) smooth $u(x,t)$ as $t \rightarrow \infty \Rightarrow$ largertime step and larger mesh size.

In general, we have the following estimation

$$\begin{cases} \|u(t)\| \leq \|u^0\| & t \in I \\ \left\| \dot{u}(t) \right\| \leq \frac{c}{t} \|u^0\| & t \in I \end{cases} \quad (3)$$

Exercise: (1) prove (3) by using (**)

(2) prove (3) using Energy method

$$\begin{aligned} u \left(\frac{du}{dt} - \frac{d^2u}{dx^2} \right) &= 0 \stackrel{(9)}{\Rightarrow} \int \frac{1}{2} \frac{du^2}{dt} + \int |\nabla u|^2 = 0 \\ &\Rightarrow \frac{1}{2} \frac{d}{dt} \|u\|^2 = -\|\nabla u\|^2 \stackrel{\text{poincare ineq.}}{<} -\|u\|^2 \\ &\Rightarrow \|u\| < e^{-t} \|u^0\| \Rightarrow \|u\| < \|u^0\| \end{aligned}$$

Exercise: using Energy method to prove $\left\| \dot{u} \right\| \leq \frac{1}{t} \|u^0\|$

Semi-discretization in space for (+) with dirichlet data ($\Gamma_2 = \phi$)

$$\text{consider } u_h(x, t) = \sum_{i=1}^m \psi_i(t) \varphi_i(x) \quad \left(\begin{array}{l} \text{in previous example} \\ \psi_i(t) = c_w e^{-w^2 t} \\ \varphi_i(x) = \sin wx \quad m = \infty \end{array} \right)$$

$$\begin{aligned} &\Rightarrow \sum_{i=1}^m \dot{\psi}_i(x) \underbrace{(\varphi_i, \varphi_j)}_M + \sum_{i=1}^m \psi_i(t) \underbrace{\langle \nabla \varphi_i, \nabla \varphi_j \rangle}_K \\ &\quad \left\{ \begin{array}{l} = \underbrace{\langle f(t), \varphi_j \rangle}_F, \quad j = 1 \sim m \\ \sum_{i=1}^m \psi_i(0) \underbrace{(\varphi_i, \varphi_j)}_{U^0} = \underbrace{(u^0, \varphi_j)}_{U^0} \quad j = 1 \sim m \end{array} \right. \quad \vec{\psi} = (\psi_1, \psi_2, \dots, \psi_m) \\ &\stackrel{(4)}{\Rightarrow} \begin{cases} M \dot{\vec{\psi}} + K \vec{\psi} = F(t) \\ M \vec{\psi}(0) = U^0 \end{cases} \Rightarrow \text{system of ordinary equation} \end{aligned}$$

Recall that the condition of M ($\chi(M) = \Delta(1)$)

and the condition of K ($\chi(k) = o(h^{-2})$) as $h \rightarrow 0$

consider $\vec{w} = e^{-i\vec{w}t} \vec{v}$ for the homogeneous case $M \vec{w} + K \vec{w} = 0$

$$\Rightarrow \left[-(Mi) \vec{w} + K \right] \vec{v} = 0 \quad \Rightarrow \quad \text{the eigenvalue } w \text{ and } V_w$$

Solve the
generalized
eigenvalue
problem gives

the general solution can now be expressed as

$$\vec{\psi} = \sum_{j=1}^m c_j e^{i w_j t} v_j \quad \text{here } \vec{c} = \{c_j\}_{j=1 \sim m} \text{ satisfies } M \cdot V \cdot \vec{c} = U^0, \quad V = [v_1, v_2, \dots, v_m]$$

Similar to the model problem,

one has smooth mode corresponding to the $\min_{j=1 \sim m} \{w_j\} \approx o(1)$

and highly oscillatory mode corresponding to the $\max_{j=1 \sim m} \{w_j\} \approx o(h^{-2})$

$\vec{\psi}$ has components live on time scales from $\underset{\text{(smooth mode)}}{o(h^{-2})}$ to $\underset{\text{(oscillatory mode)}}{o(1)}$

(or say from $o(1)$ to $o(h^2)$)

As mentioned previously, one should use methods which adapt the size of each time step according to the smoothness of $\vec{\psi}$. Moreover, if explicit time discretization is employed, one should choose the step

as small as $o(h^2)$ in order to prevent potential instability. $\left(\begin{array}{l} \Rightarrow \text{truncation error} \\ \Rightarrow o(\Delta t) = o(h^2) \end{array} \right)$

To avoid extreme small time step, implicit time discretization such as implicit Euler or Crank-Nicolson method can be employed.