

Numerical methods for the heat equation

$$\left\{ \begin{array}{l} \text{PDE: } \partial_t w(x,t) - \Delta w(x,t) = f(x,t) \quad (x,t) \in \Omega \times (0,T) \\ \text{b.c.: } w(x,t)|_{\partial\Omega} = g(x,t) \quad " \in \Gamma \times (0,T) \\ \text{i.c.: } w(x,0) = u_0(x) \quad x \in \Omega \end{array} \right.$$

We interpret $u(x,t)$ as a function of t which maps a time t to some function $\tilde{u}(x)$, $t \mapsto \tilde{u}(x)$ $\tilde{u}(t)(x) = u(x,t)$. $u(t, \cdot) = u(t)$

$$\text{PDE: } \frac{d}{dt} w(t) + Au(t) = f(t) \Rightarrow \text{ODE - setting where } -\Delta = A$$

$$\Rightarrow \frac{d}{dt} u(t) = \tilde{f}(t, u(t)) = f(t) - Au(t).$$

• Try linear numerical methods for ODEs now.

• Implicit Euler: $\left\{ \begin{array}{l} \tilde{f} = \frac{I}{\tau} \quad \tau \in \mathcal{N}, \quad t_w = w \cdot \tau. \\ \text{for } w = 1, 2, \dots \quad ; \text{ compute } \\ \frac{w(t_{n+1}) - w(t_n)}{\tau} = \tilde{f}(t_{n+1}, u(t_{n+1})) = f(t_{n+1}) - Au(t_{n+1}). \end{array} \right.$

• Combine with FE \mathcal{A}

For each t_n , we consider the PDE:

$$u(x, t_{n+1}) - \mathcal{S} \Delta u(x, t_{n+1}) = \mathcal{S} f(x, t_{n+1}) + \mathcal{S} u(x, t_n)$$

$$\textcircled{*}_1 \begin{cases} (\mathcal{I}d - \mathcal{S} \Delta) u(x, t_{n+1}) = \dots \\ \text{b.c. } u(x, t_{n+1}) = \underbrace{g(x, t_{n+1})}_{V_{g^{n+1}}} \text{ on } \partial \Omega \end{cases}$$

$\textcircled{*}_2$ Weak formulation $\Rightarrow (u(t_n), \sigma)_{\Omega}$ Find $u(t^{n+1}) \in \underbrace{H_{g^{n+1}}^1(\Omega)}_{V_{g^{n+1}}}$ s.t. $\forall \sigma \in \underbrace{H_0^1(\Omega)}_{V_0}$ it holds

$$\int_{\Omega} u(t_{n+1}) \sigma \, dx + \mathcal{S} \underbrace{(\nabla u(t_{n+1}), \nabla \sigma)_{\Omega}}_{a(u, \sigma)} = \mathcal{S} \underbrace{(f^{n+1}, \sigma)_{\Omega}}_{V_0} + (u(t^n), \sigma)_{\Omega}$$

Discrete weak formulation

$\textcircled{*}_3$ Use FE \mathcal{A} to solve $\textcircled{*}_2$, leading to

Find $u_n(t_{n+1}) \in V_{n, g^{n+1}}$ s.t. $\forall \sigma_n \in V_{n, 0}$ it holds that

$$(u_n(t_{n+1}), \sigma_n)_{\Omega} + \mathcal{S} a(u_n(t_{n+1}), \sigma_n) = \dots$$

④ linear algebra formulation.

$$\{\varphi_i\}_{i=1}^W = \text{dim } V_W \delta^{n+1}$$

{ to simplify notation, $g = 0$ }

$$u_w(t_{n+1}) = \sum_{j=1}^W U_j^{n+1} \varphi_j(x)$$

$$u_w(t_n) = \sum_{j=1}^W U_j^n \varphi_j$$

• For $i = 1 \dots W$:

$$\underbrace{\sum_{j=1}^W U_j^{n+1} \underbrace{(\varphi_j, \varphi_i)}_{\mu_{ij}}}_{(u_w^{n+1}, \varphi_i)_\Omega} + \sigma \sum_{j=1}^W U_j^{n+1} \underbrace{a(\varphi_j, \varphi_i)}_{=: A_{ij}} = \underbrace{\sigma \sum_{j=1}^W U_j^{n+1} \varphi_j}_{b_i^{n+1}} + \underbrace{\sum_{j=1}^W U_j^n \underbrace{(\varphi_j, \varphi_i)}_{\mu_{ij}}}_{(u_w^n, \varphi_i)_\Omega}$$

$$U^{n+1} := (U_j^{n+1})_{j=1}^W, \quad U^n :=$$

• Initial conditions. $\int_{\Omega} u_0(x) = u_w^0(x) = \sum_{j=1}^W U_j^0 \varphi_j(x) \quad U^0 = (U_j^0)_{j=1}^W.$

For $n = 1, 2 \dots$ \mathcal{R}_i solve

$$\mathcal{M} U^{n+1} + \sigma A U^{n+1} = \sigma b^{n+1} + \mathcal{R} U^n.$$

General theta method

- covers backward Euler, forward Euler and Crank-Nicolson method.

Given w^n and $\tau_n = t^{n+1} - t^n$ (for simplicity we assume a fixed time step $\tau = \tau_n \forall n$),

find w^{n+1} s.t.

$$M \frac{(w_{n+1} - w_n)}{\tau} + \theta A w^{n+1} + (1-\theta) A w^n = \theta f^{n+1} + (1-\theta) f^n$$

- Recall that θ -methods can be interpreted as trying to approximate w at $t_n + \theta \tau$, so the rhs of 1) could be replaced by $f^{n+\theta} = f(t_n + \theta \tau)$.

θ -methods reduces to

$\theta = 0$ explicit/forward Euler:

Advantages / Disadvantages:

first order in time, cheap to compute but not A-stable, leads to severe time-step restrictions

$\theta = \frac{1}{2}$ Crank-Nicolson:

second order in time, implicit, A-stable but not strongly A-stable / not very dissipative

very popular scheme

$\theta = 1$ implicit/backward Euler

first order, strongly A-stable, but very dissipative, not well-suited to compute stationary flows.

A quiz glimpse at numerical methods for parabolic OCP

• State equation (SE)

$$\partial_t y - \nabla \cdot (\mathcal{X} \nabla y) = \beta \overset{\text{distributed control}}{\downarrow} w \quad \text{in } \Omega \times (0, T) =: Q \quad \text{space-time cylinder}$$

$$-\mathcal{X} \partial_n y = 0 \quad \text{on } \partial\Omega \times (0, T)$$

$$y(0) = y_0 \quad \text{on } \Omega \times \{0\}$$

• Cost functional

$$J(y, w) = \underbrace{\frac{1}{2} \int_0^T \int_{\Omega} (y - y_d^e)^2 dx dt}_{= \|y - y_d^e\|_Q^2} + \underbrace{\frac{\alpha_1}{2} \int_{\Omega} (y(T) - y_d^e)^2 dx}_{\text{end-time observations}} + \underbrace{\frac{\alpha_2}{2} \int_0^T \int_{\Omega} w^2 dx dt}_{\|w\|_Q^2}$$

• Example of a parabolic OCP: minimize $J(y, w)$ subject to SE and $w \in \mathcal{U}_{ad} = \mathcal{L}^2(Q)$.

• Adjoint equation (AE)

$$-\partial_t p - \nabla \cdot (\mathcal{X} \nabla p) = \bar{y} - y_d^e \quad \text{in } Q$$

$$-\mathcal{X} \partial_n p = 0 \quad \text{on } (0, T) \times \partial\Omega$$

$$p(T) = \alpha_1 (\bar{y}(T) - y_d^e)$$

$$\left(\frac{p(T), \varphi_i}{\mathcal{X} \bar{p}} \right) = \alpha_1 (\bar{y}(T) - y_d^e, \varphi_i) = \alpha_1 (\bar{y} - y_d^e, \varphi_i)$$

• Optimality conditions

$$(\partial_2 \bar{w} + \beta p, w - \bar{w})_Q \geq 0$$

$$\int_0^T \int_{\Omega} (\partial_2 \bar{w} + \beta p)(w - \bar{w}) dx dt \geq 0 \quad \forall w \in \mathcal{U}_{ad}$$

• Riesz representation of $J'(w)$ ($J(w) = J(y(w), w)$)

$$\nabla J(w) = \alpha_2 w + \alpha p$$

$$(y_d^e, \varphi_i) = (\pi y_d^e, \varphi_i)$$

• Optimize the discretize (ODE) approach (for $\delta_1 = 0$, no end-time observation, unconstrained w)

OC: $(\delta_2 \bar{w} + \beta p, w)_Q = 0 \Leftrightarrow (\delta_2 \bar{w}(t) + \beta p(t), w(t))_Q = 0$ for a.e. t in $(0, T)$.

$\Leftrightarrow \delta_2 \bar{w} + \beta p = 0$ for a.e. (x, t) in Q .

• We discretize (SE), (AE), (OC) using FEM and then backward Euler.

• FEM: $V_w \subseteq V = H^1(\Omega)$, $\{\varphi_i\}_{i=1}^{N_v}$ basis functions

$U_w \subseteq W = L^2(\Omega)$ $\{\varphi_i\}_{i=1}^{N_w}$. Simplify by assuming $U_w = V_w$. $\varphi_i = \psi_i$

$y_w(x, t) = \sum_{j=1}^{N_v} y_j(t) \varphi_j(x)$, $u_w(x, t) = \sum_{j=1}^{N_v} u_j(t) \varphi_j(x)$, $\vec{y}(t) = (y_i(t))_{i=1}^{N_v}$, \vec{w}

• SE (discrete) $(\partial_t y_w, \varphi_i)_Q + (k \nabla y_w, \nabla \varphi_i)_Q = (\beta u_w, \varphi_i)$ $i=1, \dots, N_v$.

• SE (algebraic) $M \dot{\vec{y}} + A \vec{y} = \beta M \vec{w}$ $\vec{y}(0) = (\vec{y}_i^0)_{i=1}^{N_v}$

$M_{ij} = (\varphi_j, \varphi_i)_Q$ $A_{ij} = (k \nabla \varphi_j, \nabla \varphi_i)_Q$ $V_w \ni \vec{w} = \sum_{j=1}^{N_v} w_j^0 \varphi_j$

• AE (algebraic) $-M \dot{\vec{p}} + A^T \vec{p} = M (\vec{y} - \vec{y}_a^e)$ $M \vec{p}(T) = M (\delta \vec{y} - \vec{y}_a^e)$

coefficients from an interpolation/projection of y_a^e

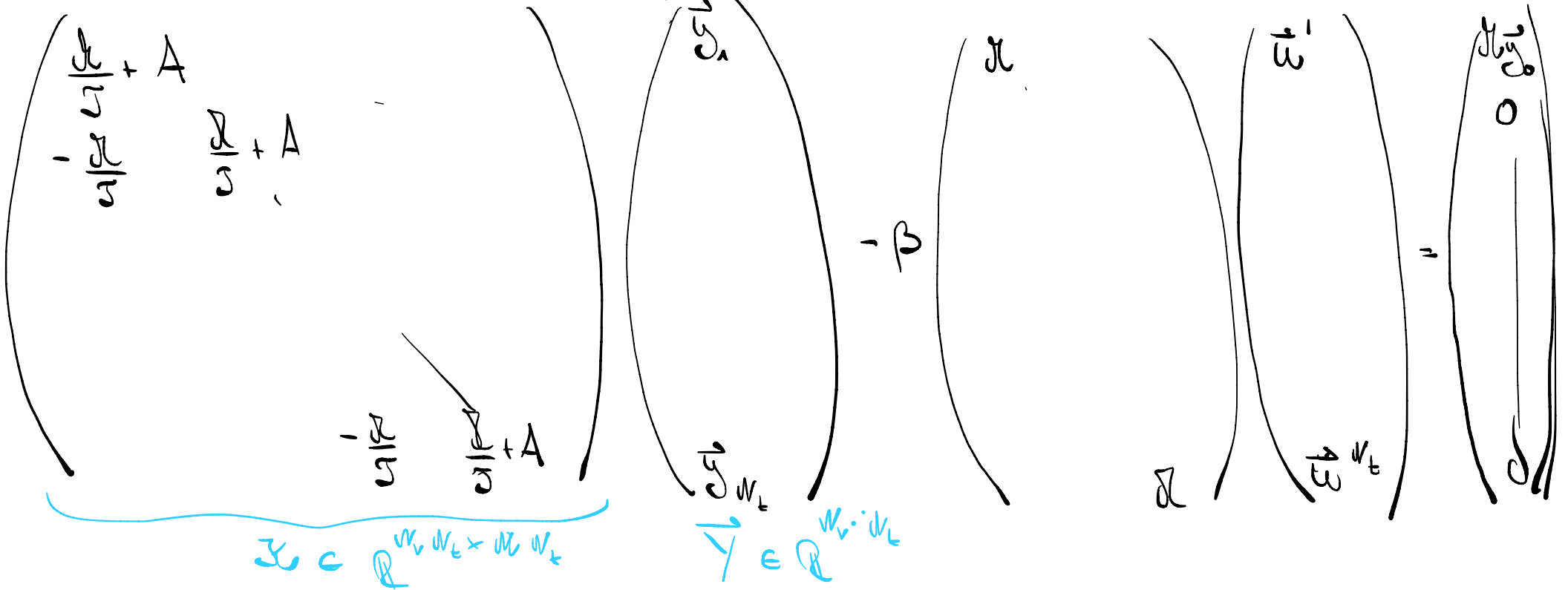
• OC (algebraic) $\delta_2 M \vec{w} + \beta A^T \vec{p} = 0$ for $t \in (0, T)$ \vec{y}_a^e collects dof-coeff from L^2 -projection.

• Discrte in time using Backward Euler: $\sigma = \frac{T}{s}$

SE $\mathcal{L} \bar{y}^{n+1} + \sigma A \bar{y}^{n+1} = \sigma \beta \mathcal{L} \bar{w}^{n+1} + H \bar{y}^n$ $w = 0, \dots, N_t$

$$\left(\frac{\mathcal{L}}{\sigma} + A \right) \bar{y}^{n+1} - \frac{\mathcal{L}}{\sigma} \bar{y}^n - \beta \mathcal{L} \bar{w}^{n+1} = 0$$

• Collect everything in big matrix system:



• How large is this system. $A, \mathcal{L} \in \mathbb{R}^{N_t \times N_t}$

Same discretization (BE) for AE: $-\mathcal{L} \dot{p} + A \dot{p} = M(\dot{y} - \dot{y}_a)$

Time-discretization $\frac{\mathcal{L}}{\Delta t}(\dot{p}^n - \dot{p}^{n+1}) + A^T \dot{p}^n = \mathcal{L}(\dot{y}^n - (\dot{y}_a)^n)$

$n = 0, 1, \dots, N-1, 0.$

$$\begin{pmatrix} \frac{\mathcal{L}}{\Delta t} + A^T & & & \\ & \ddots & & \\ & & \frac{\mathcal{L}}{\Delta t} + A^T & \\ & & & \frac{\mathcal{L}}{\Delta t} + A^T \end{pmatrix} \begin{pmatrix} \dot{p}^0 \\ \vdots \\ \dot{p}^{N-1} \end{pmatrix} + \begin{pmatrix} \mathcal{L} \\ \vdots \\ \mathcal{L} \end{pmatrix} = \begin{pmatrix} \mathcal{L} \dot{y}^0 \\ \vdots \\ \mathcal{L} \dot{y}^{N-1} \end{pmatrix}$$

$$= \begin{pmatrix} \mathcal{L}(\dot{y}_a)^0 \\ \vdots \\ \mathcal{L}(\dot{y}_a)^{N-1} \end{pmatrix}$$

3 $w_y \cdot w_t$

OC:

$$\begin{pmatrix} \alpha & & & \\ & \ddots & & \\ & & \alpha & \\ & & & \alpha \end{pmatrix} \begin{pmatrix} \dot{y}^1 \\ \vdots \\ \dot{y}^{N_t} \end{pmatrix} + \beta \begin{pmatrix} \mathcal{L}^T \\ \vdots \\ \mathcal{L}^T \end{pmatrix} \begin{pmatrix} \dot{p}^1 \\ \vdots \\ \dot{p}^{N_t} \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

Projected gradient methods

- Minimize reduced functional $f(w) := J(y(w), w)$ using the projected gradient method as before

Algorithm

- Initial control u_1

- For $n = 1, 2, \dots$

- Compute state y_n solving the forward heat equation
- adjoint state p_n solving the backward equations

- Compute new descent direction $d_n = -\nabla f(u_n)$ via Riesz representatio

$$d_n = -(\delta_2 u_n + \beta p_n)$$

$$(\delta_2 \bar{u} + \beta p, u - \bar{u})_Q \geq 0$$

- Find appropriate step length α_n by trying to solve

$$\Rightarrow (\delta_2 \bar{u}(x,t) + \beta p(x,t))(u(x,t) - \bar{u}(x,t))$$

$$f(P_{u_{\text{ad}}}(u_n + \alpha_n d_n)) = \min_{\alpha > 0} f(P_{u_{\text{ad}}}(u_n + \alpha d_n)).$$

$$\geq 0 \text{ for a.e. } (x,t) \in Q.$$

e.g. using a suitable backtracking algorithm.

- Note: as before for (time-dependent) box-constraints $\bar{f}_0(x,t) \leq w(x,t) \leq \bar{f}_1(x,t)$, OC translates to

$$\bar{w}(x,t) = P_{[\bar{f}_0(x,t), \bar{f}_1(x,t)]} \left(-\frac{\beta}{\delta_2} p(x,t) \right).$$