

Tidsstegning 1

Implicit (backward) Euler: Har betraktat ekvationen

$$\dot{u} - \nabla \cdot a \nabla u = f,$$

med $a = a(x)$ och $f = f(x, t)$, som efter multiplikation med $v = v(x)$ och integration i tid och rum och partiell integration i rumsled kunde skrivas

$$\int_{\Omega} v \underbrace{\int_{I_n} \dot{u}}_{u_n - u_{n-1}} + \int_{\Omega} \nabla v \cdot a \nabla \underbrace{\int_{I_n} u}_{\approx k_n u_n} = \int_{\Omega} v \underbrace{\int_{I_n} f}_{\approx k_n f_n},$$

där $u_n = u(x, t_n)$ och $f_n = f(x, t_n)$.

.. p.1/38

Tidsstegning 3

För enkelhets skull betraktar vi fortsättningsvis motsvarande **skalära** ekvationen

$$\dot{u} + a u = f,$$

där a är en given konstant och $f = f(t)$. Motsvarande tidsdiskretiseringsmetod för beräkning av $U_n \approx u_n = u(t_n)$ blir

$$U_n + k a U_n = U_{n-1} + k f_n,$$

dvs

$$U_n = (1 + k a)^{-1} (U_{n-1} + k f_n),$$

där $f_n = f(t_n)$.

.. p.3/38

Tidsstegning 1

Implicit (backward) Euler: Har betraktat ekvationen

$$\dot{u} - \nabla \cdot a \nabla u = f,$$

med $a = a(x)$ och $f = f(x, t)$, som efter multiplikation med $v = v(x)$ och integration i tid och rum och partiell integration i rumsled kunde skrivas

$$\int_{\Omega} v \underbrace{\int_{I_n} \dot{u}}_{u_n - u_{n-1}} + \int_{\Omega} \nabla v \cdot a \nabla \underbrace{\int_{I_n} u}_{\approx k_n u_n} = \int_{\Omega} v \underbrace{\int_{I_n} f}_{\approx k_n f_n},$$

där $u_n = u(x, t_n)$ och $f_n = f(x, t_n)$.

.. p.1/38

Tidsstegning 2

Med $u_n \rightarrow U_n = \sum_{j=1}^m U_{n,j} \phi_j$ och $v = \phi_i$ gav detta diskretiseringsmetoden

$$M (U_n - U_{n-1}) + k_n A U_n = k_n F_n,$$

där alltså $U_n \approx u_n$, $k_n = t_n - t_{n-1}$ är tidssteget, M är massmatrisen, A styvhets/diffusionsmatrisen, och F_n är lastvektorn med element $\int_{\Omega} \phi_i f_n$.

Efter multiplikation med M^{-1} kan detta skrivas

$$U_n + k \tilde{a} U_n = U_{n-1} + k \tilde{f},$$

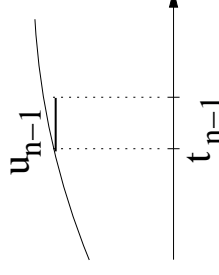
där $k = k_n$, $\tilde{a} = M^{-1} A$, och $\tilde{f} = M^{-1} F_n$.

.. p.2/38

Tidsstegning 4

Explicit (forward) Euler: Vi erinrar oss också att motsvarande approximationer

$$\int_{I_n} \dot{u} \underbrace{\approx k u_{n-1}}_{u_n - u_{n-1}} + \int_{I_n} u \underbrace{\approx k f_{n-1}}_{\approx k f_{n-1}} = \int_{I_n} f,$$



.. p.4/38

Tidsstegning 5

leder till diskretiseringsmetoden

$$U_n - U_{n-1} + k a U_{n-1} = k f_{n-1},$$

dvs

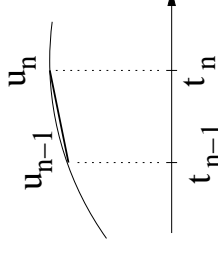
$$U_n = (1 - k a) U_{n-1} + k f_{n-1},$$

allmänt känd som Explicit (forward) Euler.

-- P.6/38

Tidsstegning 7

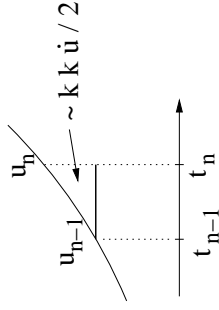
Crank-Nicolson: Uppenbarligen finns möjligheten att approximera integralerna $\int_{I_n} u$ och $\int_{I_n} f$ bättre än ovan med t.ex. *trapezregeln*:



-- P.7/38

Tidsstegning 6

De båda Euler metoderna kan förväntas vara jämförbara vad gäller *noggrannhet*. De fel vi introducerar genom att ersätta $\int_{I_n} u$ och $\int_{I_n} f$ med $k u_n$ och $k f_n$ resp $k u_{n-1}$ och $k f_{n-1}$ är ju i båda fallen av storleksordning k^2 , dvs felet efter *en tidsenhet*, dvs efter k^{-1} tidssteg, kan väntas ha accumulerats till storleksordning k .



-- P.6/38

Tidsstegning 8

dvs

$$\int_{I_n} \dot{u} + a \int_{I_n} u = \int_{I_n} f, \quad \int_{I_n} u \approx k (u_{n-1} + u_n) / 2 \approx k (f_{n-1} + f_n) / 2$$

vilket leder till diskretiseringsmetoden

$$U_n - U_{n-1} + k a (U_{n-1} + U_n) / 2 = k (f_{n-1} + f_n) / 2,$$

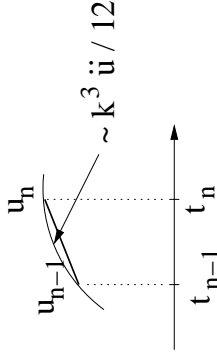
dvs

$$U_n = (1 + \frac{k}{2} a)^{-1} (U_{n-1} - \frac{k}{2} a U_{n-1} + k (f_{n-1} + f_n) / 2).$$

-- P.8/38

Tidsstegning 9

Felet för denna metod bör vara av storleksordning k^3 per tidssteg, dvs av storleksordning k^2 efter en tidsenhet.



.. p.10/38

Tidsstegning 11

Motsvarande metoder för värmeledningsproblemet ovan är

$$\begin{aligned}(M + kA)U_n &= U_{n-1} \\ MU_n &= (M - kA)U_{n-1} \\ (M + \frac{k}{2}A)U_n &= (M - \frac{k}{2}A)U_{n-1}\end{aligned}$$

Vi noterar att alla dessa metoder kräver **ekvationslösning**, men om massmatrisen M i vänsterledet ersätts med motsvarande **lumpade massmatris** \bar{M} så kan metod II skrivas

$$U_n = \bar{M}^{-1}(M - kA)U_{n-1},$$

dvs metoden blir **explicit**, dvs man erhåller utan ekvationslösning en **formel** för U_n , därav namnet **explicit Euler**.

.. p.11/38

Tidsstegning 10

För $f = 0$ reduceras de tre metoderna till:

$$\begin{aligned}U_n &= (1 + ka)^{-1}U_{n-1} \\ \bar{U}_n &= (1 - ka)U_{n-1} \\ U_n &= (1 + \frac{k}{2}a)^{-1}(1 - \frac{k}{2}a)U_{n-1}\end{aligned}$$

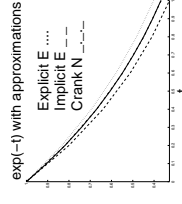
.. p.10/38

Tidsstegning 12

Egenskaper: Inleder med att studera modellproblemet

$$\dot{u} + u = 0 \quad t > 0, \quad u(0) = 1,$$

med exakt lösning $u(t) = \exp(-t) = 1 - t + \frac{1}{2}t^2 - \frac{1}{6}t^3 + \dots$



Man noterar inga väsentliga skillnader mellan Euler metoderna, medan Crank-Nicolsson som väntat är mera exakt.

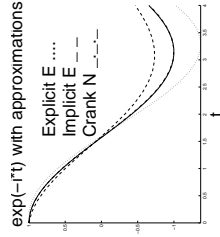
.. p.12/38

Tidsstegning 13

Skillnaden mellan metoderna blir mera uppenbar om vi betraktar ett par andra **typ-problem**. Vi börjar med

$$\dot{u} + i u = 0 \quad t > 0, \quad u(0) = 1,$$

med (komplexvärd) lösning $u(t) = \cos(t) + i \sin(t)$, med realdel och imaginärdel motsvarande lägeskoordinat resp hastighet hos en odämpad svängande massa med massan 1 upphängd i en fjäder med fjäderkonstant 1.



--p.13/38

Tidsstegning 15

Vi erinrar oss att (den preliminära) felanalysen för CN indikerade ett fel av storleksordning $k^3 \dot{u}/12$ per tidssteg, dvs $0.2^3 40^2/12 \approx 1.1$ i det aktuella fallet, så resultatet borde kanske inte vara överraskande, men **oscillationen** är slående!

För implicit Euler är resultatet måhända något bättre än väntat, eftersom (den preliminära) analysen pekade på ett möjligt fel av storleksordning $k^2 \dot{u}/2 = 0.8$.

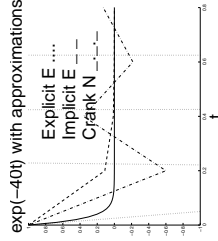
--p.15/38

Tidsstegning 14

Slutligen betraktar vi ett **styvt** problem

$$\dot{u} + 40 u = 0 \quad t > 0, \quad u(0) = 1,$$

med lösning $u(t) = \exp(-40 t)$.



--p.14/38

Vi noterar att Crank-Nicolsson nu uppför sig konstigt. Implicit Euler är bättre! Explicit Euler är helt ute!

Tidsstegning 16

Vi kan förstå uppförandet hos de tre metoderna genom att studera ett enskilt tidssteg av längd 0.2.

$$U_n = (1 + 0.2 \cdot 40)^{-1} U_{n-1} = 1/9 U_{n-1}$$

$$U_n = (1 - 0.2 \cdot 40) U_{n-1} = -7 U_{n-1}$$

$$U_n = (1 + \frac{0.2}{2} \cdot 40)^{-1} (1 - \frac{0.2}{2} \cdot 40) U_{n-1} = -3/5 U_{n-1}$$

med faktorer framför U_{n-1} som ju bör approximera $\exp(-k a) = \exp(-8)$ (eftersom för den exakta lösningen $u_n = \exp(-a t_n) = \exp(-a (t_{n-1} + k)) = \exp(-k a) \exp(-a t_{n-1}) = \exp(-k a) u_{n-1}$) men som för explicit Euler och CN blir negativa, och för Euler dessutom blir stor vilket leder till instabilitet! Problemet ligger i att kvantiteten $k a$ här inte är tillräckligt "liten".

--p.16/38

Tidsstegning 17

Notera att faktorn $(1 - k a)$ utgör de två första termerna i Taylor approximationen av $\exp(-k a)$, som ger en okey approximation om (och endast om i detta fall) $k a$ är "liten". Analogt utgör förstas faktorn $(1 + k a)^{-1} = 1 - k a + k^2 a^2 - k^3 a^3 + \dots$ motsvarande implicit Euler en okey approximation till $\exp(-k a)$ om $k a$ är liten. Notera att $\exp(-k a) = 1 - k a + \frac{1}{2} k^2 a^2 - \dots$. En analys visar att motsvarande faktor för CN sammanfaller med denna serie upp t.o.m. k^2 -termen, vilket än en gång visar att metoden är av ordning k^3 fel i varje tidssteg.

-- p.17/28

Tidsstegning 19

Samma typ av CN-reaktion kan man se som resultat av alla "plötsliga förändringar" i data. Notera att begynnelsevärdet $u(0) = 1$ i fallet $\dot{u} + 40 u = f$ kan ses som en stationär lösning motsvarande $f = 40$, och att sedan, vid tiden $t = 0$ lasten f plötsligt ändras till $f = 0$.

Ett sätt att undvika detta uppförande för CN är förstås att helt enkelt undvika problem med denna typ av snabba "ryck", t.ex. genom att påföra kraft/lastförändringar **successivt**.

-- p.19/28

Tidsstegning 18

Vi studerar nu motsvarigheten till CN's uppförande på ett värmeledningsproblem.

-- p.18/28

Tidsstegning 20

En intressant kombination av Implicit Euler's goda dämpande egenskaper och CN högre noggrannhet är att inleda med **två** Implicita Eulersteg, för att därefter fortsätta med CN.

-- p.20/28

Tidsstegning 21

Andra ordningens tidsderivator: Vi betraktar modelproblemet

$$\ddot{u} + a u = 0 \quad t > 0, \quad u(0) = u_0, \quad \dot{u}(0) = v_0,$$

vilket vi skriver om som ett ekvivalent **system** med $v = \dot{u}$:

$$\begin{cases} \dot{u} - v = 0 \\ \dot{v} + a u = 0 \end{cases}$$

$$\text{dvs med } w = [uw]^\top \text{ och } A = \begin{bmatrix} 0 & -1 \\ a & 0 \end{bmatrix}$$

$$\dot{w} + A w = 0 \quad t > 0, \quad w(0) = \begin{bmatrix} u_0 \\ v_0 \end{bmatrix}.$$

--p.21/38

Tidsstegning 23

Galerkin metoder:

Galerkinmetoder bygger på att utifrån en given **lösningsansats**, dvs "försökslösningar" \tilde{U} av viss typ, "testa" sig fram till en lösning U vars **residual** uppfyller lämpliga ortogonalitetsvillkor.

För vår modellekvation $\dot{u} + a u = f$ ges residualen av $\dot{U} + a U - f$ och ortogonaliteten kommer till uttryck som

$$\int_{I_n} v(\dot{U} + a U - f) = 0, \quad \text{dvs} \quad \int_{I_n} v(\dot{U} + a U) = \int_{I_n} v f,$$

för alla v av viss typ.

--p.23/38

Andra ordningens tidsderivator: Vi betraktar modelproblemet

$$\ddot{u} + a u = 0 \quad t > 0, \quad u(0) = u_0, \quad \dot{u}(0) = v_0,$$

vilket vi skriver om som ett ekvivalent **system** med $v = \dot{u}$:

$$\begin{cases} \dot{u} - v = 0 \\ \dot{v} + a u = 0 \end{cases}$$

$$\text{dvs med } w = [uw]^\top \text{ och } A = \begin{bmatrix} 0 & -1 \\ a & 0 \end{bmatrix}$$

$$\dot{w} + A w = 0 \quad t > 0, \quad w(0) = \begin{bmatrix} u_0 \\ v_0 \end{bmatrix}.$$

--p.21/38

Tidsstegning 22

Heuns metod: Explicit Euler blir här

$$W_n = W_{n-1} - A W_{n-1}.$$

Ett försök till komponentvis implementering skulle kunna se

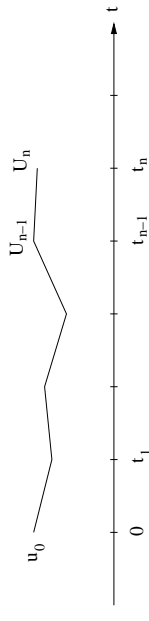
```
ut:
while time < finaltime
u=u+k*v;
v=v-k*a*u;
time=time+k;
end
```

Vad är **fel** med denna implementering?
Fungerar implementeringen? Jämför med en korrekt implementering.

--p.22/38

Tidsstegning 24

cG1: Här står c för **continuous**, G för **Galerkin**, och 1 för grad 1, eller **linjär**, dvs utgångspunkten är en **ansats** till lösning på I_n som i figur, med just dessa egenskaper:



och U_n bestäms så att

$$\int_{I_n} \dot{U} + a U = \int_{I_n} f,$$

motsvarande residualortogonalitet mot alla **konstanta** $v = v(t)$ på I_n .

--p.24/38

Tidsstegning 25

Speciellt erhåll för a konstant

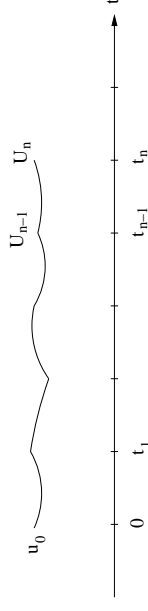
$$\underbrace{\int_{I_n} \ddot{U}}_{U_n - U_{n-1}} + \underbrace{\int_{I_n} aU}_{a k_n (U_{n-1} + U_n)/2} = \int_{I_n} f,$$

dvs metoden sammanfaller väsentligen med CN (medelvärde av f i högerledet motsvaras i CN av trapezapproximationen $k_n (f_{n-1} + f_n)/2$). För f linjär i t är metoderna identiska.

..p.26/38

Tidsstegning 27

cG2: Analogt utgår cG2 från en styckvis **kvadratisk** ansats:



och U på I_n bestäms av att

$$\int_{I_n} v(\ddot{U} + aU) = \int_{I_n} v f,$$

för $v = \psi_{n-1}$ och $v = \psi_n$, vilket tillsammans med kontinuitetskravet ger tre ekvationer för det kvadratiska polynomet U på I_n .

..p.27/38

Tidsstegning 26

Mera allmänt kan ansatsen U på I_n representeras med hjälp av tidsbasfunktionerna $\psi_{n-1}(t) = (t_n - t)/k_n$ och $\psi_n(t) = (t - t_{n-1})/k_n$ som $U(t) = U_{n-1}\psi_{n-1}(t) + U_n\psi_n(t)$, vilket ger

$$U_n - U_{n-1} + \underbrace{\int_{I_n} a\psi_{n-1}U_{n-1}}_{\frac{k_n}{2}\tilde{a}_{n-1}} + \underbrace{\int_{I_n} a\psi_nU_n}_{\frac{k_n}{2}\tilde{a}_n} = \int_{I_n} f,$$

vilket för a konstant ju reduceras till

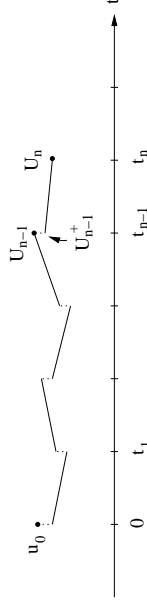
$$U_n - U_{n-1} + \frac{k_n}{2}U_{n-1} + \frac{k_n}{2}U_n = \int_{I_n} f.$$

..p.26/38

Tidsstegning 28

dG1: Om vi släpper kontinuitetskravet erhålls s.k. dG metoder, där d för **discontinuous**, och G för Galerkin, som förut.

För dG1 använder vi en ansats som i figur, där vi på det nya intervallet I_n har att bestämma både U_{n-1}^+ och U_n .



..p.28/38

Tidsstegning 29

För detta behöver vi två ekvationer och vi väljer

$$\int_{I_n} v(\dot{U} + aU) = \int_{I_n} v f,$$

med $v = \psi_{n-1}$ och $v = \psi_n$ som förut. Notera emellertid att nu ingår i \dot{U} på I_n även derivatan av det initiala hoppet $U_{n-1}^+ - U_{n-1}$, vilket ger bidraget

$$(U_{n-1}^+ - U_{n-1})v(t_{n-1})$$

till integralen. Innan vi räknar på ett konkret dG1-exempel noterar vi att även styckvis konstant ansats är möjligt utan kravet på kontinuitet:

-- p.30/38

Tidsstegning 31

För att bestämma $U_n = U(t)$ på I_n kräver vi nu att

$$\int_{I_n} v(\dot{U} + aU) = \int_{I_n} v f,$$

för $v = v(t) = 1$, dvs

$$(U_{n-1}^+ - U_{n-1}) + \int_{I_n} (\underbrace{\dot{U}}_{=0} + aU) = \int_{I_n} f,$$

dvs

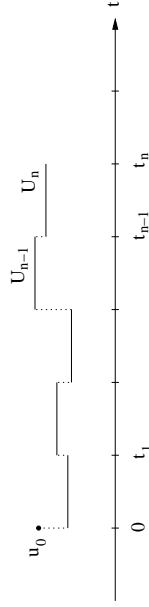
$$U_n - U_{n-1} + \int_{I_n} aU_n = \int_{I_n} f.$$

Vi noterar att för a och f konstanta sammanfaller denna metod med Implicit Euler.

-- p.31/38

Tidsstegning 30

dG0: Här utgår vi alltså från en ansats som i figur med U **styckvis konstant**, med $U_{n-1}^+ = U_n = U(t)$ för alla $t \in I_n$:



-- p.30/38

Tidsstegning 32

Vi återgår nu till dG1 med a konstant och $f = 0$. För $v = \psi_{n-1}$ erhålls

$$U_{n-1}^+ - U_{n-1} + \int_{I_n} \psi_{n-1}(U_n - U_{n-1}^+) + a \int_{I_n} \psi_{n-1}(U_{n-1}^+ \psi_{n-1} + U_n \psi_n)$$

dvs

$$U_{n-1}^+ - U_{n-1} + \frac{k}{2}(U_n - U_{n-1}^+) + a\left(\frac{k}{3}U_{n-1}^+ + \frac{k}{6}U_n\right) = 0,$$

och för $v = \psi_n$

$$\frac{k}{2}(U_n - U_{n-1}^+) + a\left(\frac{k}{6}U_{n-1}^+ + \frac{k}{3}U_n\right) = 0.$$

-- p.32/38

Tidsstegning 33

Efter elimination av U_{n-1}^+ erhålls

$$U_n = \left(1 + \frac{2ak}{3} + \frac{a^2 k^2}{6}\right)^{-1} \left(1 - \frac{ak}{3}\right) U_{n-1},$$

vilket ger oss en ny (rationell) approximation av av exponentialfaktorn $\exp(-ak)$.

-- p.33/38

Tidsstegning 35

Vi utgår från *felekvationen*

$$\int_0^T v(\dot{e} + e) = 0, \quad (e = u - U)$$

här giltig för alla $v = v(t)$ som är styckvis konstanta. Med ϕ sådan att

$$-\dot{\phi} + \phi = 0 \quad \text{for } t < T, \quad \phi(T) = e(T),$$

fås

-- p.35/38

Tidsstegning 34

A posteriori felanalys för cG1:

Vi betraktar ekvationen $\dot{u} + u = f$, $t > 0$, $u(0) = u_0$ och visar att för cG1 approximationen U till u gäller

$$|(u - U)(T)| \leq \max_{[0,T]} |k| (f - \dot{U} - U).$$

-- p.34/38

Tidsstegning 36

$$\begin{aligned} |e(T)|^2 &= \phi(T) e(T) + \int_0^T \underbrace{(-\dot{\phi} + \phi)}_{=0} e \\ &= \int_0^T \phi(\dot{e} + e) = \int_0^T (\phi - v)(\dot{e} + e) = \int_0^T (\phi - v) \underbrace{(f - \dot{U} - U)}_{=: r} \\ &\leq \int_0^T k^{-1} |\phi - v| \max_{[0,T]} |k r(U)|. \end{aligned}$$

-- p.36/38

Tidsstegning 37

Med lämpligt val av interpolant v av ϕ , känd uppskattning av interpolationsfelet $\phi - v$, och genom att utnyttja att $\dot{\phi} = \phi$ & $\phi(t) = \exp(t - T)e(T)$, erhålls

$$|e(T)|^2 \leq \underbrace{\int_0^T \exp(t - T) dt}_{=1 - \exp(-T) \leq 1} \max_{[0, T]} |kr(U)|,$$

dvs

$$|e(T)| \leq \max_{[0, T]} |kr|.$$

..p.37/38

Tidsstegning 38

Övning: Generalisera till fallet $\dot{u} + au = f$, med a konstant > 0 .

Övning: Är samma a posteriori uppskattning giltig även för ekvationen $\dot{u} - u = f$?

..p.38/38