

STEFFEN HAUG

# *Numeriske Metodar*

Prosjekt 3

# I. KVADRATURMETODAR

Vi skal implementere to forskjellige kvadraturmetodar: Ein adaptiv kvadraturmetode som tilpassar steglengda slik at ein er meir nøye med tilnærmingar over intervall som er utsatt for større feil, og Rombergs metode som burkar ekstrapoleringsteknikkar for å konvergere raskare.

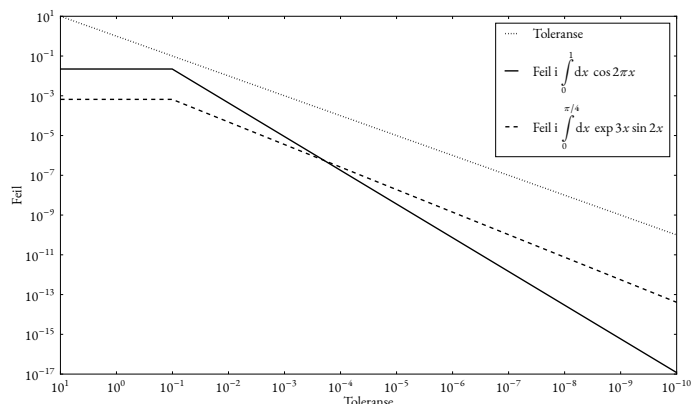
## i. ADAPTIV SIMPSON-KVADRATUR

Adaptive kvadraturmetodar baserer seg på å integrere over eit intervall med ein enkel kvadraturmetode, estimere feilen, for so å rekursivt dele intervallet i to og beregne eit meir nøyaktig estimat med halvparten av steglengda dersom feilen er uakseptabelt stor.

Algoritma (**APPENDIX A**), med Simpsons metode over kvart del-intervall, vart nytta til å estimere integrala i **TABELL I**. Ei oversikt over samanhengen mellom feil og toleranse ser ein i **FIGUR I**.

**TABELL I:** Estimasjonar v. h. a. adaptiv Simpson-kvadratur med toleranse  $\varepsilon = 10^{-10}$ .

INTEGRAL	ESTIMAT
$\int_{x=0}^{x=1} dx \cos 2\pi x$	$-2.7138785046 \times 10^{-17}$
$\int_{x=0}^{x=\pi/4} dx \exp 3x \sin 2x$	2.58862863250



**FIGUR I:** Feil i adaptiv Simpson-kvadratur. Ein ser at feilen er grensa ovanfrå av toleransen.

## ii. ROMBERG-KVADRATUR

Romberg-algoritma (**APPENDIX B**) organiserer (i prinsippet) dei utrekna verdiane i ei nedre-triangular matrise, der kvar *kolonne* er estimat av integralet med aukande presisjon (kortare steglengde) nedover langs kolonna. Trikset til Romberg er at ein – ved hjelp av to estimat av ein gitt orden, med forskjellige steglengder  $h$  og  $h/t$  – kan ein konstruere eit estimat av *høgare orden*, ved hjelp av *Richardson-ekstrapolering*. Ekstrapoleringa aksellererer konvergensens til følgja, som ein ser i **FIGUR 2**.

La  $A(h)$  vere ein approksimasjon på ekvidistante nodar med steglengde  $h$

$$A(h) \approx I := \int_a^b dx f(x)$$

slik at  $A(h) - I = \mathcal{O}(h^n)$ , der restfunksjonen er eit polynom i  $h$ . I tilfelle for Rombergs metode er denne gitt ved Euler-Maclaurin-ekspansjonen. Generelt, ved å bruke tilnærmingar med to forskjellige steglengder,  $A(h)$  og

$A(b/t)$ , har vi

$$A(b/t) - A(b) = \left( I + k \left( \frac{b}{t} \right)^n + \mathcal{O}(b^{n+1}) \right) - \left( I + k b^n + \mathcal{O}(b^{n+1}) \right)$$

her har vi skrive ut det første leddet i rest-funksjonen; målet er å bli kvitt  $b^n$ -leddet. Dette oppnår vi dersom vi gonger  $A(b/t)$  med  $t^n$ :

$$t^n A(b/t) - A(b) = t^n I - I + \mathcal{O}(b^{n+1})$$

Her har vi innlemma  $t^n$  i det asymptotiske leddet. Del med  $t^n - 1$  og trekk  $I$  frå begge sider:

$$\frac{t^n A(b/t) - A(b)}{t^n - 1} - I = \mathcal{O}(b^{n+1}) \quad (1)$$

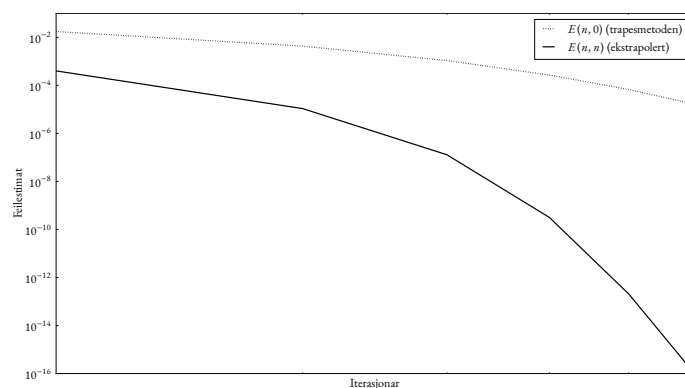
Altso har vi, ved å bruke *to* approksimasjonar med ulike steglengder, klart å hoste opp ein approksimasjon av *hø-gare orden*! Dersom vi startar med  $N$  estimat v. h. a. trapesmetoden kan vi ekstrapolere  $N$  gongar, og organisere approksimasjonane i ei nedretriangulær matrise. Organiserer vi utrekningane slik at kvar ekstrapolasjon  $R(m, n)$  kun brukar dei to approksimasjonane  $R(m-1, n)$  og  $R(m-1, n-1)$  kan vi beregne approksimasjonane rad for rad, i staden for kolonne for kolonne. På denne måten oppnår vi ønska presisjon med ferrast mulig nodar i den opprinnelige trapesmetoden, og vi treng kun lagre to rader i gangen, så vi sparar minne. Denne metoden heiter *Rombergs metode*, og nokre estimat av integral er i

**TABELL 2.**

**TABELL 2:** Estimasjonar v. h. a. Rombergs metode med toleranse  $\varepsilon = 10^{-10}$ .

INTEGRAL	ESTIMAT
$\int_{x=0}^{x=1} dx \cos 2\pi x$	$3.30152344409 \times 10^{-17}$
$\int_{x=0}^{x=1} dx \sqrt[3]{x}$	$0.74998846514 +$
$\operatorname{erf}(1) = \frac{2}{\sqrt{\pi}} \int_{x=0}^{x=1} dx e^{-x^2}$	$0.84270079294$

Det ser ikkje ut til at toleransen  $\varepsilon$  er ei øvre skranke for feilen i Rombergs metode. Feilen i (+) er om lag  $\delta = 1.15 \times 10^{-5}$ , som er større enn  $\varepsilon = 10^{-10}$ . For det same integralet, med den same toleransen, gjer den adaptive kvadraturmetoden svaret 0.75, som er eksakt, men tar om lag ti gongar so lang tid å kjøre.



**FIGUR 2:** Feil i Rombergs metode. Det er tydeleg at elementa langs diagonalen konvergerer raskt samanlikna med elementa i første kolonne. Den akselererte følgja når presisjonen etter berre 6 iterasjonar.

## II. SIMULASJON AV FRI, STIV LEKAM

Vi ønsker å simulere ein *fri, stiv lekam*, med andre ord ein lekam som ikkje lar seg deformere, som roterer fritt i rommet frå ein gitt start-tilstand.

Fri rotasjon betyr at netto påført dreiemoment er null. At lekamen ikkje lar seg deformere medfører at treighetsmomentet er konstant. (massen “flyttar seg ikkje” relativt til rotasjons-aksen) Dette betyr at normen til dreieimpulsen er bevart (den er fri til å endre retning), og rotasjonsenergien  $E = \frac{1}{2}T\omega^2$  er bevart. [1]

### i. DEFINISJON AV PROBLEMET

Differensiallikninga [1]

$$T \frac{d\omega}{dt} + \omega \times T\omega = \mathbf{M} \quad (2)$$

skildrar rotasjonen til eit stivt legeme. Her er  $\mathbf{M}$  påført dreiemoment,  $T$  treighetsmoment, og  $\omega$  vinkelfart. Per antagelse er  $\mathbf{M} = 0$ . Innfør  $\mathbf{m} = T\omega$ , som er eit produkt av eit moment,  $T$ , og ein fart,  $\omega$ . Altso er  $\mathbf{m}$  ein form for *impuls* – nemleg *dreieimpulsen*. Vi har

$$\frac{d\mathbf{m}}{dt} = T \frac{d\omega}{dt} \quad \text{og} \quad \omega = T^{-1}\mathbf{m}. \quad (3)$$

I (2), trekk frå  $d\mathbf{m}/dt$  på begge sider, og snu kryssproduktet for å få like forteikn. Samtidig innfører vi notasjonen  $\dot{f}$  for tids-deriverte.

$$\dot{\mathbf{m}} = \mathbf{m} \times T^{-1}\mathbf{m} \quad (2^*)$$

Bevaringslovene nevnt til å byrje med, saman med (3), gir

$$\gamma = \|\mathbf{m}\|^2 = \mathbf{m}(t) \cdot \mathbf{m}(t) \quad (4)$$

$$E = \frac{1}{2}\mathbf{m}(t) \cdot T^{-1}\mathbf{m}(t) \quad (5)$$

Der (4) er ei sfære med radius  $\|\mathbf{m}\|$ , og (5) er ei ellipse. Sidan både  $\gamma$  og  $E$  er konstantar er løysingane  $\mathbf{m}$  begrensa til å sitje på skjæringa mellom flatene.

Anta  $T = \text{diag}(I_1, I_2, I_3)$ . Ettersom  $T$  er diagonal er  $T^{-1} = \text{diag}(1/I_1, 1/I_2, 1/I_3)$ . Skriv  $\mathbf{m} = (x \ y \ z)$ , og skriv ut (2\*) på komponentform:

$$\begin{pmatrix} \dot{x}(t) \\ \dot{y}(t) \\ \dot{z}(t) \end{pmatrix} = \begin{pmatrix} x(t) \\ y(t) \\ z(t) \end{pmatrix} \times \begin{pmatrix} 1/I_1 & 0 & 0 \\ 0 & 1/I_2 & 0 \\ 0 & 0 & 1/I_3 \end{pmatrix} \begin{pmatrix} x(t) \\ y(t) \\ z(t) \end{pmatrix}$$

$\mathbf{m} \quad T^{-1} \quad \mathbf{m}$

Rekn ut kryssproduktet:

$$\begin{pmatrix} \dot{x}(t) \\ \dot{y}(t) \\ \dot{z}(t) \end{pmatrix} = \begin{pmatrix} Ay(t)z(t) \\ Bx(t)z(t) \\ Cx(t)y(t) \end{pmatrix} \quad (6)$$

$\mathbf{m} \quad f(t, \mathbf{m})$

der  $A, B$  og  $C$  er konstantar

$$A = 1/I_3 - 1/I_2$$

$$B = 1/I_1 - 1/I_3$$

$$C = 1/I_2 - 1/I_1$$

### ii. IMPLISITT RUNGE-KUTTA

#### MIDTPUNKT-METODE

Vi ønsker å løyse differensiallikningar av sorten

$$\dot{y} = f(t, y)$$

numerisk, ved hjelp av Runge-Kutta metoden

$$y_{n+1} = y_n + hf \left( t_n + \frac{h}{2}, \frac{1}{2}(y_n + y_{n+1}) \right), \quad (7)$$

kalt *implisitt midtpunktm metode*, fordi  $y_{n+1}$  avheng av eit estimat for  $y_{n+1/2}$  (derav implisitt) for å estimere tangenten i midtpunktet mellom  $t_n$  og  $t_{n+1}$ . Substituer

$$u = \frac{1}{2}(y_n + y_{n+1}) \implies y_{n+1} = 2u - y_n$$

for å forenkle notasjonen litt. Dette gir likningssystemet

$$2u - y_n = y_n + hf(t + h/2, u)$$

som vi kan skrive

$$u = y_n + \frac{h}{2} f(t + h/2, u) \quad (8)$$

og dette må løysast med omsyn til  $u$  for kvart tidssteg. Systemet vi er interresert i har tre variablar;  $x$ ,  $y$  og  $z$ . Med  $\mathbf{u} = (x \ y \ z)$  kan vi definere

$$\mathbf{F}(\mathbf{u}) = \mathbf{u}_n + f(t_n + h/2, \mathbf{u}) - \mathbf{u} \quad (9)$$

der  $\mathbf{u}_n$  er førre iterasjon, og søke etter nullpunkt ved hjelp av ein kjend løysar for ulineære likningssystem, til dømes Newtons metode. (som er kjend frå tidlegare øvingar, men eg gjentar det mest grunnleggande for kompletthets skyld)

Gitt Jacobian-matrissa til  $\mathbf{F}$ ,  $\mathcal{J}_{\mathbf{F}}$ , kan vi bruke den iterative prosessen

$$\mathbf{u} \leftarrow \mathbf{u} - \mathcal{J}_{\mathbf{F}}^{-1} \mathbf{F}(\mathbf{u})$$

med  $\mathbf{u} = \mathbf{u}_n$  som startverdi. Iterasjonen er slutt når antal iterasjonar har nådd ei øvre grense, eller når  $\|\mathbf{F}(\mathbf{u})\| < \varepsilon$ , for ein gitt toleranse  $\varepsilon$ . Med  $y_{n+1/2} = \mathbf{u}$  reknar vi ut  $y_{n+1}$ :

$$y_{n+1} = y_n + hf(t + h/2, y_{n+1/2}) \quad (10)$$

og dette fullfører éin iterasjon av Runge-Kutta-metoden.

### iii. IMPLEMENTASJON AV RK-METODEN

Vi brukar eit symbolsk verktøy for å rekne ut (10) i forkant, til dømes `sympy`. Den eine haka ved denne implementasjonen er at differensiallikninga vi vil løyse må vere definert berre ved hjelp av funksjonar som er kompatible

med `sympy`, men det er enklare å ha eit eksplisitt uttrykk for Jacobian-matrissa ved kvart tidssteg.

Det er enklare å behandle uttrykk der dimensjonen er kjend. `sympy` har funksjonalitet for å lage vektorar med symbol, så det er (i prinsippet) mogleg å implementere metoden for  $N$  dimensjonar, men element i desse vektorane får obskure navn som gjer feilsøking til eit mareritt. Implementasjonen (**APPENDIX C**) er derfor begrensa til tre dimensjonar.

### iv. MODIFISERT EULERS METODE

Vi implementerer òg den eksplisitte RK-metoden (**APPENDIX D**)

$$y_{n+1/2} = y_n + \frac{h}{2} f(t_n, y_n)$$

$$y_{n+1} = y_{n+1/2} + hf(t_n + 1/2, y_{n+1/2})$$

for å samanlikne den implisitte metoden med ein eksplisitt metode. Metoden treng ingen forklaring, sidan den er eksplisitt.

### v. NUMERISK LØYSING AV LIKNINGA

I **LISTING I** vert bruk av løysarane demonstrert.

**LISTING I:** *Bruk av RK-løysarane*

---

```

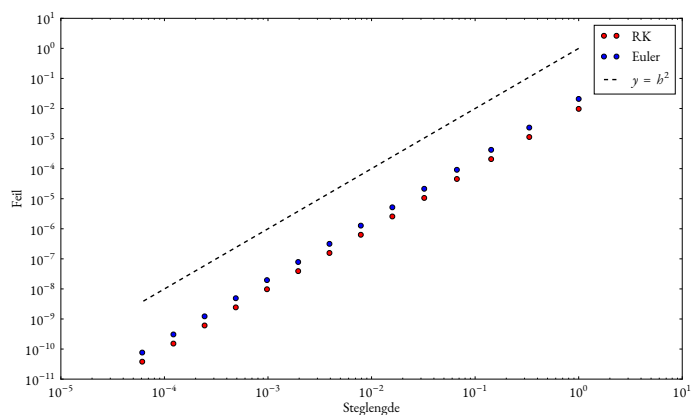
1 def f(t, m):
2     x, y, z = m
3     return [ (1/I3 - 1/I2) * y * z ,
4              (1/I1 - 1/I3) * x * z ,
5              (1/I2 - 1/I1) * x * y ]
6
7 t = np.linspace(0, 150, 150)
8
9 y0 = np.array([2, 3, 4], dtype=float)
10 y0 *= 1/linalg.norm(y0)
11
12 x, y, z = RK(f, y0 = y0, t = t)
```

---

## vi. KVADRATISK KONVERGENS

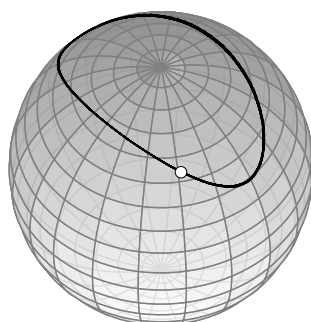
Vi brukar ei svært nøyaktig algoritme frå `sympy`, med svært kort steglengde, for å estimere  $\mathbf{m}(t = 1)$  so nøyaktig som mogleg. Dersom vi antar at denne verdien er om lag eksakt, samanlikna med verdiane vi oppnår med våre metodar med relativt lange steglengder kan vi undersøkje korleis feilen oppfører seg når steglengda minkar. Feilen kan ein sjå i **FIGUR 3**, med steglengda som abscisse.

Ved sidan er der teikna ei linje  $\gamma = h^2$ . I  $\log\log$ -plottet har denne stigningstal  $m = 2$ . ( $\log x^2 = 2 \log x$ ) Det er heilt klart at begge metodana har same stigningstal som linja, altså er begge kvadratisk konvergente i  $h$ , RK med litt lågare konstantfaktor enn Euler.

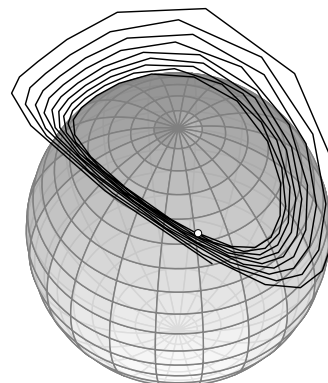


**FIGUR 3:** Feilen i begge metodane er tydeleg kvadratiske.

Merk at dette tidsrommet er svært kort. Vi er òg interesserte i å undersøkje kvar som skjer over lange tidsrom.



**(A)** Løysing v. h. a. implisitt midtpunkt-metode. Løysinga er presis, sjølv etter svært lang tid.



**(B)** Løysing v. h. a. eksplisitt modifisert Eulers metode. Løysinga vert ubrukeleg nesten med ein gang.

**FIGUR 4:** Løysingar av  $(2^*)$  med initialverdier  $\mathbf{T} = \text{diag}(1, 2, 5)$  og  $\mathbf{m}_0 = (2 \ 3 \ 4)$ . Tidsrommet er  $t = 150$ , med  $N = 150$  steg. Initialverdien  $\mathbf{m}_0$  er markert. Løysingane er plotta på ei sfære  $\gamma = \|\mathbf{m}_0\|^2$ .

## vii. LANGE TIDSROM

Løysingane produsert av løysarane er plotta i **FIGUR 4A** (implisitt RK) og **FIGUR 4B** (eksplisitt modifisert Euler). Det er tydeleg at den implisitte metoden produserer langt meir presise løysingar til dette initialverdiproblemet over lang tid: Den eksplisitte løysinga forlét sfæra nesten umiddelbart, og feilen veks raskt til å verte urimeleg stor.

## REFERANSAR

<sup>1</sup>J. R. Lien og G. Løvhøiden, *Generell fysikk for universiteter og høyskoler* (2015).

## A. SIMPSON-KVADRATUR

LISTING 2: *Adaptiv Simpson-kvadratur*


---

```

1  def S(f, a, b):
2      return (f(a) + 4 * f((a + b) / 2) + f(b)) \
3          * (b - a) / 6
4
5  def asm_quad(f, a, b, tol=1e-5):
6      I0 = S(f, a, b)
7      c = (a + b) / 2
8
9      I = S(f, a, c) + S(f, c, b)
10     err = abs(I - I0) / 15
11
12     if err < tol:
13         return I + (I - I0) / 15
14     else:
15         return asm_quad(f, a, c, tol=tol/2) \
16             + asm_quad(f, c, b, tol=tol/2)

```

---

## B. ROMBERG-KVADRATUR

LISTING 3: *Romberg-kvadratur*


---

```

1  def romberg(f, a, b, MAX_ITER=100, tol=1e-5):
2      R = np.zeros((2, MAX_ITER))
3      Rp, Rn = 0, 1
4
5      h = b - a
6      R[Rp, 0] = 0.5 * h * (f(a) + f(b))
7
8      for n in range(1, MAX_ITER):
9          h = h * 0.5
10         L = np.linspace(a + h, b - h, 1 << n - 1)
11         R[Rn, 0] = R[Rp, 0] / 2 + np.sum(f(L)) * h
12
13         for k in range(1, n + 1):
14             E = (R[Rn, k - 1] - R[Rp, k - 1]) \
15                 / ((1 << 2 * k) - 1)
16             R[Rn, k] = R[Rn, k - 1] + E
17
18         Rp, Rn = Rn, Rp
19
20         if abs(E) < tol: break
21
22     return R[Rp, n]

```

---



## C. IMPLISITT RUNGE-KUTTA

sympy-koden for likningssystemet er ikkje spesielt vanskelig, men den er lang og stygg, så eg har klipt den ut. Algoritma kan finnast i jupyter-notebook: `rk3d.pynb`, og der er sjølvsagt sympy-koden inkludert.

**LISTING 4:** *Implisitt Runge-Kutta midtpunktmetode.*

---

```

1  def RK(f, y0, t):
2
3      # (Sympy-junk removed here...)
4
5      y = np.empty((len(t), 3))
6      y[0] = y0
7
8      for n in range(len(t) - 1):
9          h = t[n + 1] - t[n]
10
11         u = newton(
12             lambda u: F(u, y[n], t[n] + h/2, h),
13             lambda u: Ji(u, y[n], t[n] + h/2, h),
14             y[n]
15         )
16
17         m = np.array(f(t[n] + h/2, u), dtype=float)
18
19         y[n + 1] = y[n] + h * m
20
21     return y.T

```

---

I metoden inngår Newtons metode; denne har vi implementert i ei tidlegare øving, men legg den ved for kompletthets skuld:

**LISTING 5:** *Newton's metode frå tidlegare øving.*

---

```

1  def newton(F, Ji, u0, tol=1e-10, maxiter=10):
2      u = u0
3      for _ in range(maxiter):
4          u = u - Ji(u).dot(F(u))
5
6          if linalg.norm(F(u)) < tol:
7              break
8
9      return u

```

---

## D. MODIFISERT EULER

LISTING 6: *Modifisert Eulers metode (eksplisitt)*

---

```
1 def euler(f, y0, t):
2     y = np.empty((len(t), 3))
3     y[0] = y0
4
5     for n in range(len(t) - 1):
6         h = t[n + 1] - t[n]
7
8         u = y[n] + np.array(f(t[n], y[n])) * h / 2
9         m = np.array(f(t[n] + h/2, u), dtype=float)
10        y[n + 1] = y[n] + h * m
11
12    return y.T
```

---