

Kap. 6 i kompendiet: Simulering av differensligninger på datamaskin og avrundingsfeil

Definisjon 6.7

En likning på formen

$$(*) \quad X_{n+k} = g(n) + f_0(n)X_n + f_1(n)X_{n+1} + \dots + f_{k-1}(n)X_{n+k-1}$$

sies å være linear, inhomogen, k-ordens differensligning

Kan simuleres på datamaskin, som for 1. orden, 2. orden, når man har initialbetingelser x_0, x_1, \dots, x_{k-1}

For oss opp til nå: $f_0(n), f_1(n) \dots$ har vært konstanter.

Hva kan gå galt i simuleringer som (*)?

1. Koeffisientene i (*) kan ikke representeres eksakt.

2. Initialbetingelsene i (*) kan ikke representeres eksakt.

Eksempel med avrundingsfeil

$$\text{Vi ser på } X_{n+2} - \frac{19}{3}X_{n+1} + 2X_n = 0, \quad X_0 = 1, \quad X_1 = \frac{1}{3}.$$

Hva skjer på datamaskin? Hva er løsningen?

Løsning:

Analytisk løsning: Kar. ligning: $r^2 - \frac{19}{3}r + 2 = 0$

$$r = \frac{\frac{19}{3} \pm \sqrt{\frac{361}{9} - 8}}{2}$$

$$= \frac{\frac{19}{3} \pm \frac{17}{3}}{2} = \frac{19 \pm 17}{6}$$

$$\Rightarrow r_1 = 6, \quad r_2 = \frac{1}{3}$$

$$\underline{\text{Generell løsning}} \quad X_n = C 3^{-n} + D 6^n$$

Init. betingelser:

$$\begin{aligned} x_0 &= 1 : \quad C + D = 1 \\ x_1 &= \frac{1}{3} : \quad \frac{1}{3}C + 6D = \frac{1}{3} \end{aligned} \quad \left. \begin{array}{l} C + D = 1 \\ \frac{1}{3}C + 6D = \frac{1}{3} \end{array} \right\} \quad \left. \begin{array}{l} D = 0, C = 1 \\ C + 18D = 1 \end{array} \right\}$$

$$\Rightarrow \underline{x_n = 3^{-n}}$$

Hvorfor? Løsningen er ustabil i den forstand at hvis vi endrer litt på initialbetingelsene, så vil den nye løsningen gå mot ∞ :

Her $x_0 = 1$ siden $\frac{1}{3}$ ikke kan rep. eksakt med 64-bits flyttall.
 $x_1 = \frac{1}{3} + \epsilon$ $(\epsilon \approx 10^{-16} \text{ for 64-bits flyttall})$

Løsning med nye initial betingelser:

$$\tilde{X}_n = C 3^{-n} + D 6^n, \text{ der } CD \text{ løser}$$

$$\begin{aligned} C + D &= 1 & D &= 1 - C \\ \frac{1}{3}C + 6D &= \frac{1}{3} + \epsilon & \frac{1}{3}C + 6(1 - C) &= \frac{1}{3} + \epsilon \Rightarrow \frac{17}{3}C = \frac{17}{3} - \epsilon \\ & & \Rightarrow C &= 1 - \frac{3}{17}\epsilon \\ & & D &= \frac{3}{17}\epsilon \end{aligned}$$

$$\text{Def følger at } \tilde{X}_n = \left(1 - \frac{3}{17}\epsilon\right)3^{-n} + \frac{3}{17}\epsilon 6^n$$

For stor n vil $\tilde{X}_n \approx \frac{3}{17}\epsilon 6^n$, slik at $x_n \rightarrow \infty$.

Anta x_n er første som blir ∞ (int). Neste blir

$$x_{n+1} = \underbrace{\frac{19}{3}x_n}_{\infty} - 2x_{n-1} = \infty$$

$$\text{Neste: } x_{n+2} = \frac{19}{3}x_{n+1} - 2x_n, \text{ som er } \infty - \infty = \text{nan.}$$

Merk Avrundingsfeil ; koeffisientene har samme effekt som over : Selv om x_0, x_1 kan representeres eksakt, så vil x_2 da få avrundingsfeil , som vil forplantet seg slik som vi så avrundingsfeilen for x_1 i eksemplet over gjorde.

Noen observasjoner:

1. Anta $r^2 + br + c = 0$ har to forskjellige reelle røtter r_1 og r_2 , og at $|r_1| > |r_2|$

Hvis \tilde{x}_n er den numeriske løsningen, så vil $\tilde{x}_n \approx Cr_1^n$ for store n , høst det skjer avrundingsfeil.

2. Vi har at

$$\frac{\tilde{x}_n}{x_{n-1}} \approx \frac{Cr_1^n + Dr_2^n}{Cr_1^{n-1} + Dr_2^{n-1}} \approx \frac{Cr_1^n}{Cr_1^{n-1}} = r_1.$$

Dette kan være utgangspunkt for å finne røtter i polynomer.

Eksempel (tilsynelatende uten avrundingsfeil)

$$2x_{n+2} - 5x_{n+1} + 2x_n = -1 \quad x_0 = 2, x_1 = \frac{3}{2}$$

Hva skjer når vi simulerer dette på datamaskin?

Løsning: Analytisk løsning :

partikulær løsning: prøv $x_n^p = A$:

$$2A - 5A + 2A = -1$$

$$-A = -1$$

$$A = 1 \Rightarrow x_n^p = 1$$

Homogen løsning: $2r^2 - 5r + 2 = 0$
 $\Rightarrow r = \frac{5 \pm \sqrt{25 - 16}}{4} = \frac{5 \pm 3}{4} \Rightarrow r_1 = 2, r_2 = \frac{1}{2}$

$$\Rightarrow x_n^h = C 2^n + D 2^{-n}$$

Generell løsning: $x_n = x_n^p + x_n^h = 1 + C 2^n + D 2^{-n}$

Spesiell løsning: $x_0 = 2 \quad \left\{ \begin{array}{l} 1 + C + D = 2 \\ C + D = 1 \end{array} \right.$
 $x_1 = 3/2 \quad \left\{ \begin{array}{l} 1 + 2C + \frac{1}{2}D = \frac{3}{2} \\ 2C + \frac{1}{2}D = \frac{1}{2} \end{array} \right. \quad \underbrace{\qquad}_{4C + D = 1}$

$$\left. \begin{array}{l} C + D = 1 \\ 4C + D = 1 \end{array} \right\} \quad C = 0, D = 1 \quad \Rightarrow x_n = 1 + 2^{-n}$$

Numerisk løsning: Maskinen regner ut $x_{n+2} = \frac{5}{2}x_{n+1} - x_n - \frac{1}{2}$

Her er det ikke avrundingsfeil, hverken i

koeffisentene, eller initialverdiene.

Så maskinen skal regne riktig, til å begynne med.

Vi får likevel avrundingsfeil til slutt.

Til slutt får vi overflow , deretter overflow ,
 -Inf -Inf
 til slutt nan.

Oppgave 18, midtveis 2021

Vi ser på $x_{n+1} - 3x_n = 1$, $x_1 = 1$

Hva skjer når vi simulerer dette med 64 bits flyttall?

Løsning: numerisk løsning: Maskinen regner ut $x_{n+1} = 3x_n + 1$.
 Ingen avrundingsfeil i starten, vi vil få overflow
 unsett.

Analytisk løsning: $x_n^p = A \Rightarrow A - 3A = 1$
 $-2A = 1$
 $A = -\frac{1}{2}$

$$r - 3 = 0 \quad x_n^h = C 3^n$$

Generell løsning: $X_n = -\frac{1}{2} + C3^n$

$$X_1 = 1 : \quad 1 = -\frac{1}{2} + 3C \Rightarrow 3C = \frac{3}{2} \Rightarrow C = \frac{1}{2}$$

$$X_n = -\frac{1}{2} + \frac{1}{2}3^n = \frac{1}{2}(3^n - 1) \rightarrow \infty$$

Oppgave 19 midtveis 2021

Vi ser på $qX_{n+2} - 3X_{n+1} - 2X_n = 0 \quad X_0 = 1, X_1 = -\frac{1}{3}$

Hva skjer når vi simulerer med 64 bits flyttall?

Løsning: Analytisk løsning:

$$qr^2 - 3r - 2 = 0 \Rightarrow \dots \Rightarrow r_1 = \frac{2}{3}, r_2 = -\frac{1}{3}$$

$$\Rightarrow X_n = C\left(\frac{2}{3}\right)^n + D\left(-\frac{1}{3}\right)^n$$

$$\begin{array}{l} X_0 = 1 \\ X_1 = -\frac{1}{3} \end{array} \quad \left. \begin{array}{l} C + D = 1 \\ \frac{2}{3}C - \frac{1}{3}D = -\frac{1}{3} \end{array} \right\} \quad \left. \begin{array}{l} C + D = 1 \\ -2C + D = 1 \end{array} \right\} \quad C = 0, D = 1$$

$$\Rightarrow X_n = \underline{\left(-\frac{1}{3}\right)^n}$$

Numerisk: Maskinen regner ut $X_{n+2} = \frac{1}{3}X_{n+1} + \frac{2}{q}X_n$

Avirundingsfeil i både koeffisienter og initialverdi

Derfor vil maskinen regne ut noe på formen

$$(1-\varepsilon)\left(-\frac{1}{3}\right)^n + \varepsilon\left(\frac{2}{3}\right)^n$$

i stedet. Vil unsett gå mot 0.