

Fel och datoraritmetik

Avrundningsfel och flyttalssystem

Från labben:

- diskretiseringsfel, avrundningsfel
- Beräkningen $A^{-1}A$ blev inte riktigt enhetsmatrisen
- olika sätt att mäta fel (rel fel, absolut fel)
- begrepp ϵ_M , Inf, NaN, overflow, underflow

Avrundningsfel och flyttalssystem

Exempel)

Uttryck	blir exakt	i MATLAB
$\cos(\pi/2)$	0	6.1232e-017
$0.42 - 0.5 + 0.08$	0	-1.3878e-017
$0.08 + 0.42 - 0.5$	0	0
$A^{-1} \cdot A$	Enhetsmatrisen	Se lab

Varför blir det inte alltid exakt rätt svar när man arbetar i MATLAB?
Testa gärna även olika beräkningar på din miniräknare.

Avrundningsfel och flyttalssystem

Först hur man mäter fel:

Det exakta talet betecknas x och \hat{x} samma tal som innehåller något fel t ex fel från avrundningar eller mätfel.

- Absolut fel $|x - \hat{x}|$
- Relativt fel $\frac{|x - \hat{x}|}{|x|}$

Om x är en vektor blir det istället

- Absolut fel $\|x - \hat{x}\|$
- Relativt fel $\frac{\|x - \hat{x}\|}{\|x\|}$

$\|\cdot\|$ kallas för norm. Mer om detta senare.

Avrundningsfel och flyttalssystem

Exempel (fr labben)

- Du köper en varmkorv en lördagkväll. Den kostar 15 kr, men av misstag ger du 20 kr.

$$\text{Absolut fel: } |15 - 20| = 5$$

$$\text{Relativt fel: } \frac{|15-20|}{|15|} \approx 0.333 = 33.3\%$$

- Du köper en ny bil för 299 995 kr. Du betalar 300 000 och bryr dig inte om växeln tillbaka.

$$\text{Absolut fel: } |299995 - 300000| = 5$$

$$\text{Relativt fel: } \frac{|299995-300000|}{|299995|} \approx 0.0000166 \approx 0.0017\%$$

Ofta uppfattar man det relativa felet mer korrekt (inte alltid). I båda fallen förlorar man lika mycket, men det känns sannolikt mindre i det andra fallet.

Avrundningsfel och flyttalssystem

Representation av tal i dator:

Alla tal lagras i ett begränsat antal bitar i minnet, vanligen i binär form:

0	1	0	1	1	1	0	0
---	---	---	---	---	---	---	---

betyder $(01011100)_2 =$

$$= 0 \cdot 2^7 + 1 \cdot 2^6 + 0 \cdot 2^5 + 1 \cdot 2^4 + 1 \cdot 2^3 + 1 \cdot 2^2 + 0 \cdot 2^1 + 0 \cdot 2^0 = (92)_{10}$$

- Heltal

Inga problem. Heltal upp till en viss storlek lagras exakt.

- Reella tal

Reella tal kan inte lagras exakt utan måste avrundas eller huggas av. Representationen av reella tal i datorn kallas *flyttalsrepresentation* och talen kallas *flyttal*.

Avrundningsfel och flyttalssystem

Reella tal:

Reellt tal x kan skrivas (exakt)

$$x = m \cdot \beta^e$$

där β är den bas som används, e exponent och mantissan

$$m = \pm(0.d_1d_2\dots) = d_1\beta^{-1} + d_2\beta^{-2} + \dots, 0 \leq d_i < \beta$$

där

$$\beta^{-1} \leq |m| < 1 \Rightarrow d_1 \neq 0 \text{ kallas normaliserad form}$$

Innebär alltså att decimalpunkten flyttas så att den första termen efter decimalpunkten, d_1 blir skild från noll.

Avrundningsfel och flyttalssystem

Exempel I bas $\beta = 10$:

$$\begin{aligned} \pi &= 3.14159265\dots = 0.314159265\dots \cdot 10^1 = \\ &= (3 \cdot 10^{-1} + 1 \cdot 10^{-2} + 4 \cdot 10^{-3} + \dots) \cdot 10^1. \end{aligned}$$

Stämmer att $0 \leq d_i < 10$ och $0.1 \leq |m| < 1$.

När reella talet x lagras i datorn görs det som flyttal, betecknas $fl(x)$. Hur ska reella talet representeras i datorns begränsade antal bitar? Exponenten e kan bara lagras exakt upp till en viss storlek, och mantissan m måste kapas eftersom den är oändlig.

Avrundningsfel och flyttalssystem

Ett flyttal $fl(x)$ representeras

$$fl(x) = \hat{m} \cdot \beta^e, \quad \hat{m} = \pm(.d_1d_2, \dots, d_p), \quad \text{där} \\ 0 \leq d_i < \beta, \quad d_1 \neq 0, \quad L \leq e \leq U, \quad (0 \text{ hanteras separat})$$

Detta medför att

$|fl(x)| \geq min$, annars "underflow"

$$min = (1 \cdot \beta^{-1} + 0 \cdot \beta^{-2} + \dots + 0 \cdot \beta^{-p})\beta^L = \beta^{L-1}$$

$|fl(x)| \leq max$, annars "overflow".

$$max = ((\beta - 1)\beta^{-1} + (\beta - 1)\beta^{-2} + \dots + (\beta - 1)\beta^{-p})\beta^U = \\ (1 + \beta^{-1} + \beta^{-2} + \dots + \beta^{-p-1} - \beta^{-1} - \dots - \beta^{-p-1} - \beta^{-p})\beta^U = \\ (1 - \beta^{-p})\beta^U,$$

Avrundningsfel och flyttalssystem

- Flyttalssystemet karakteriseras alltså av (β, p, L, U) på en given dator.
- p kallas precision
- Vi antar att exponenten e lagras exakt (eftersom heltal)
- Mantissan rundas av
- \hat{m} och e lagras. β är fixt för en given dator och lagras ej. Vanligen gäller att $\beta = 2$.

Avrundningsfel och flyttalssystem

Exempel)

Flyttalssystemet $(\beta, p, L, U) = (2, 3, 0, 2)$.

\hat{m} kan anta följande positiva värden (motsvarande för negativa tal):

$$(0.100) = \frac{4}{8} \quad (0.101) = \frac{5}{8} \quad (0.110) = \frac{6}{8} \quad (0.111) = \frac{7}{8}$$

och följdaktligen fås

$$\begin{array}{l|l} e = 0 & \begin{array}{l} min = \frac{1}{2} \quad \frac{5}{8} = 0.625 \quad \frac{6}{8} = 0.75 \quad \frac{7}{8} = 0.875 \\ \frac{2}{2} = 1 \quad \frac{10}{8} = 1.25 \quad \frac{12}{8} = 1.5 \quad \frac{14}{8} = 1.75 \\ \frac{16}{8} = 2 \quad \frac{20}{8} = 2.5 \quad \frac{24}{8} = 3 \quad max = \frac{28}{8} = 3.5 \end{array} \end{array}$$

OBS! Flyttalen är *ej jämnt fördelade*! Större tal – glesare representation.

Avrundningsfel och flyttalssystem

Exempel, forts)

Givet ett tal x som ska representeras i datorn som flyttal. Hur stort kan felet maximalt bli?

Vid avrundning fås maximalt fel i mantissan då m ligger mitt mellan två tal \hat{m} i flyttalssystemet, t ex $m = \frac{11}{16}$ i exemplet. Antag avrundning uppåt. Vi får felet:

$$|\hat{m} - m| = \left| \frac{12}{16} - \frac{11}{16} \right| = \frac{1}{16}, \quad \text{dvs } |\hat{m} - m| \leq \frac{1}{16}$$

Allmänt gäller för felet i mantissan

$$|\hat{m} - m| \leq \frac{1}{2}\beta^{-p}$$

Avrundningsfel och flyttalssystem

För felet i hela talet, dvs $fl(x)$ gäller:

- Absoluta felet

$$|fl(x) - x| = |\hat{m}\beta^e - m\beta^e| = |(\hat{m} - m)\beta^e| \leq \frac{1}{2}\beta^{-p} \cdot \beta^e$$

Felet beror av storleken på talet x , en följd av att flyttalen är glesare representerade för stora tal

- Relativa felet

$$\frac{|fl(x) - x|}{|x|} \leq \frac{\frac{1}{2}\beta^{-p} \cdot \beta^e}{|m| \cdot \beta^e} = \frac{1}{2} \frac{\beta^{-p}}{|m|} \leq \frac{1}{2} \frac{\beta^{-p}}{\beta^{-1}} = \frac{1}{2} \beta^{1-p}$$

Beror *ej* på talets storlek !

Avrundningsfel och flyttalssystem

Vi har

$$\frac{|fl(x) - x|}{|x|} \leq \epsilon_M, \quad \text{där } \epsilon_M = \frac{1}{2}\beta^{1-p}$$

vid avrundning.

ϵ_M kallas *maskinepsilon* eller *avrundningsenheten* och är en maskinberoende konstant.

Det gäller att

$$\epsilon_M, \text{ är det minsta tal } \epsilon \text{ s.a. } fl(1 + \epsilon) > 1$$

I MATLAB: Kommandot \gg eps ger maskinepsilon

Avrundningsfel och flyttalssystem

Exempel

Antag (β, p, L, U) och $\beta = 2$, dvs bas 2.

Vanligt antalet bitar i minnet, p är

$$p = 24 \Rightarrow \epsilon_M \approx 0.00000006 = 6 \cdot 10^{-8}$$

$$p = 53 \Rightarrow \epsilon_M \approx 2 \cdot 10^{-16}$$

Normalt används 53 bitar, t ex av MATLAB. Detta ger som bäst ca 16 decimalers noggrannhet. Vid beräkningar av olika typ så gäller att alla relativa fel som är som är mindre än detta enbart är avrundningsfel!

Detta förklarar resultaten i slide 2!

Avrundningsfel och flyttalssystem

Vi kan nu skriva

$$\begin{aligned} \frac{|fl(x) - x|}{|x|} \leq \epsilon_M &\Rightarrow \frac{fl(x) - x}{x} = \delta, \quad |\delta| \leq \epsilon_M \\ &\Rightarrow fl(x) = x + \delta x = x(1 + \delta) \end{aligned}$$

Motsvarande uttryck för felet vid beräkningar med de fyra räknesätten:

$$\begin{aligned} \frac{|fl(x \text{ op } y) - (x \text{ op } y)|}{|x \text{ op } y|} &\leq \epsilon_M \\ \Rightarrow fl(x \text{ op } y) &= (x \text{ op } y)(1 + \delta_1), \quad |\delta_1| \leq \epsilon_M \end{aligned}$$

där *op* betecknar +, -, ·, /, dvs de fyra räknesätten.

Avrundningsfel och flyttalssystem

Avrundningsfelen vid beräkningar är vanligen mycket små i förhållande till andra fel (t ex i indata eller diskretisering).

Men, kan vara problem vid subtraktion $fl(x - y)$:

$$fl(fl(x) - fl(y)) = fl(x(1 + \delta_1) - y(1 + \delta_2)) = fl(x - y + x\delta_1 - y\delta_2) = fl((x - y)(1 + \frac{x}{x-y}\delta_1 - \frac{y}{x-y}\delta_2))$$

där $|\delta_i| \leq \epsilon_M$.

Om $x - y$ litet relativt x resp y fås stort fel ! Detta kallas för *kancellation* och uppstår då två nästan lika stora tal subtraheras. Kan ofta undvikas genom omskrivning.

Ex) $\sqrt{1+x} - \sqrt{1-x} = \frac{2x}{\sqrt{1+x} + \sqrt{1-x}}$.

Avrundningsfel och flyttalssystem

Några konsekvenser:

- ej meningsfullt med tester av typen **if x = y then**
istället **if abs(x-y) < tolerans then**
eller **if abs(x-y) < tolerans · abs(x) then**
- undvik subtrahera nästan lika stora tal
- associativa och distributiva lagarna gäller ej flyttal
Har i praktiken ingen betydelse.
- summera om möjligt termer i växande storleksordning
Sällan något man bryr sig om – har liten betydelse.

IEEE flyttalsrepresentation

- under 60 och 70-talen hade varje dator tillverkare sitt eget flyttalssystem
- flyttalsstandard utvecklades under tidigt 80-tal och följdes av tillverkare som Intel o Motorola
- utvecklat av arbetsgrupp hos the Institute for Electrical and Electronics Engineers (IEEE)
- IEEE-standarden har tre viktiga krav
 - konsistent flyttalsrepresentation
 - korrekt avrundningsaritmetik
 - konsistent behandling av exceptionella situationer

IEEE flyttalsrepresentation

- tre standardtyper: **single precision, double precision** och **extended precision**
 - **Exempel) IEEE enkel precision**
 $\boxed{\pm e_1 e_2 e_3 \dots e_8 \quad d_1 d_2 \dots d_{23}}$
tecknet 1 bit, exponenten 8 bitar, mantissan 23 bitar
 - **Exempel IEEE dubbel precision**
 $\boxed{\pm e_1 e_2 e_3 \dots e_{11} \quad d_1 d_2 \dots d_{52}}$
tecknet 1 bit, exponenten 11 bitar, mantissan 52 bitar

IEEE flyttalsrepresentation

- Maskinepsilon blir

$$\text{IEEE single} \quad \epsilon_M = 2^{-23} \approx 1.2 \cdot 10^{-7}$$

$$\text{IEEE double} \quad \epsilon_M = 2^{-52} \approx 1.1 \cdot 10^{-16}$$

$$\text{IEEE extended} \quad \epsilon_M = 2^{-63} \approx 1.1 \cdot 10^{-19}$$

- (r, p, L, U) i IEEE-standard

$$\text{IEEE single} \quad (2, 24, -126, 128)$$

$$\text{IEEE double} \quad (2, 53, -1022, 1024)$$

Obs "hidden bit"

IEEE flyttalsrepresentation

IEEE definierar 5 olika "exceptions":

- "Invalid operation", t ex $\frac{0}{0}$, $0 * \infty$, ges värdet NaN (Not a Number)
- division med 0, sätt till $\pm\infty$
- overflow, sätt till $\pm\infty$ eller största flyttal
- underflow, sätt till 0 (eller "subnormalt" tal)
- korrekt avrundning av reellt tal (ej exceptionell situation egentligen)

Diskretiseringsfel

Förutom avrundningsfel finns även diskretiseringsfel. **Exempel**
Numerisk derivering med *differenskvot*:

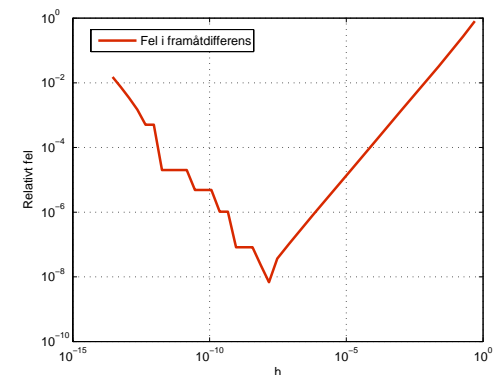
$$f'(x) \approx \frac{f(x+h) - f(x)}{h}$$

När h blir mindre borde approximationen bli successivt bättre eftersom finare indelning – mindre diskretiseringsfel.

Blir det så? Se resultat på labben.

Diskretiseringsfel

Hur blir det i praktiken?



Varför ser grafen ut som den gör?

Diskretiseringsfel och avrundningsfel

Tolkning:

- För stora h spelar diskretiseringsfelet den dominerande rollen, kan bortse från avrundningsfelet
- För små h är avrundningsfelet dominerande
- avrundningsfelet i det här fallet beror på
 - kancellation i täljaren för små h
 - division med litet tal förstärker felet i täljaren
- Avrundningsfelet uppför sig "kaotiskt", medan diskretiseringsfelet ger en jämn och snygg kurva

Diskretiseringsfel och avrundningsfel

Sammanfattning:

- Diskretiseringsfelet spelar vanligen den dominerande rollen. Det är vanligt att man helt bortser från avrundningsfelet
- Avrundningsfel får konsekvenser i vissa fall, t ex kancellation. Om möjligt bör man ta hänsyn till detta i koden genom omskrivning.
- Exakt noll existerar i princip inte i praktiken! Om två reella tal x och y anses lika, skiljer de sig ändå i någon decimal.
- Ett relativt fel i storleksordningen ϵ_M efter beräkning med flyttal är enbart slumpmässigt "skräp".