Computerphysik Programmiertutorial 4a

Prof. Dr. Matteo Rizzi und Dr. Markus Schmitt - Institut für Theoretische Physik, Universität zu Köln

Github: https://github.com/markusschmitt/compphys2022

Inhalt dieses Notebooks: Rechnen auf dem Rechner: Maschinengenauigkeit, Vergleichen von Fließkommazahlen

Rechnen auf dem Rechner: Maschinengenauigkeit Zahlen werden im Computer in einem Binärcode dargestellt und für jede Zahl steht nur eine begrenzte Anzahl von Bits zur Verfügung. Es können daher

weder alle ganzen Zahlen Z noch alle reellen Zahlen R dargestellt werden.

In [1]:

In [5]:

In [6]:

start=2^63-4

Ganze Zahlen - Int

Wir haben schon in einem früheren Tutorial gesehen, dass Ganzzahlen in 64 bits als Binärzahlen dargestellt werden. Das ergibt automatisch eine Grenze für die größte darstellbare Zahl. Schauen wir uns diese Grenze an:

 $zahl1 = 2^63$ bitstring(zahl1)

Out[1]:

In [2]: $zahl2 = 2^63-1$ bitstring(zahl2)

Out[2]:

In [3]: zahl1

-9223372036854775808 Out[3]: In [4]: zahl2

Out[4]: 9223372036854775807

Schreiben wir eine Funktion, die die n-te Harmonische Zahl H_n berechnet:

Die ganzen Zahlen auf dem Computer sind ein "Kreis":

for i **in** 1:6 println("\$start + \$i = \$(start+i)")

Fließkommazahlen - Float

function H forward(n, mytype=Float32)

S = mytype(0.0)for k in 1:n

println(H_forward(1000))

println(H backward(1000))

println(H_forward(100000)-H_backward(100000))

return S

end

7.4854784

7.4854717

0.0006980896

eps = 1.0

zu kodieren.

end

Dezimal

Dezimal

Dezimal

using Plots

xlabel!("n")

ylabel!("Differenz")

end

function H approx(n) n = Float64(n)

n werte = [2^n for n in 1:20]

Hn fwd = [H forward(2^n) for n in 1:20] $Hn_bwd = [H_backward(2^n)$ for n in 1:20]Hn approx = $[H approx(2^n)$ for n in 1:20]

using Printf

In [14]:

In [17]:

In [20]:

In [21]:

Out[21]:

In [22]:

In [23]:

In [24]:

1.32

1.2+0.12

while 1.0+eps > 1.0 eps /= 2.0

println("n = \$n")

 $println("eps = 1/2^{n} = eps")$ println("1 + eps = \$(1+eps)")println("1 + 2eps = \$(1+2eps)")

n=0

In [12]:

In [13]:

9223372036854775804 + 1 = 92233720368547758059223372036854775804 + 2 = 92233720368547758069223372036854775804 + 3 = 92233720368547758079223372036854775804 + 4 = -92233720368547758089223372036854775804 + 5 = -92233720368547758079223372036854775804 + 6 = -9223372036854775806

 $H_n = \sum\limits_{k=1}^n rac{1}{k}$

Auf dem Computer können wir sehr leicht große Summen ausrechnen. Ein Beispiel ist die Harmonische Reihe

S += mytype(1.0)/kend

 $r = \pm m \times b^e$ Das Vorzeichen wird in einem Bit kodiert, für Mantisse und Exponent steht jeweils eine feste Zahl weiterer Bits zur Verfügung. Das Kodieren der Mantisse in

den Exponenten bedeutet, dass es wie bei Ganzzahlen auch eine größte und kleinste darstellbare Fließkommazahl gibt.

und **Exponent** e zerlegt werden. Eine reelle Zahl $r \in \mathrm{R}$ wird also geschrieben als

Mit unterschiedlicher Reihenfolge der Summation erhalten wir unterschiedliche Ergebnisse!

Da der Computer nur mit einer bestimmten Zahl von signifikanten Ziffern rechnet, ist der Unterschied zwischen Zahlen nur begrenzt auflösbar. Diese "Auflösung" können wir experimentell bestimmen, indem wir fragen was die kleinste Zahl ϵ ist, so dass auf dem Computer $1.0+\epsilon>1.0$:

einer begrenzten Anzahl von Bits bedeutet, dass wir bei jeder Zahl nur eine feste Anzahl von signifikanten Ziffern kennen. Die begrenzte Anzahl von Bits für

Reelle Zahlen werden im Computer als **Fließkommazahlen** behandelt. Das bedeutet, dass sie bezüglich einer festen **Basis** b in **Vorzeichen** \pm , **Mantisse** m

n += 1 end

```
n = 53
eps = 1/2^53 = 1.1102230246251565e-16
1 + eps = 1.0
Die 64 bits des Datentyps Float64 sind wie folgt aufgeteilt (Bild gestohlen von benjaminjurke.com):
?
Wir haben also 1 Bit für das Vorzeichen, 11 Bits kodieren den Exponenten als ganze Zahl zwischen -1022 und 1023. Als Basis wird b=2 verwendet. Die
darstellbaren Zahlen bewegen sich also (in etwa) zwischen 2^{-1022}pprox 10^{-308} und 2^{1023}pprox 10^{308}. Die übrigen 52 Bits werden verwendet um die Mantisse als
```

 $m=1+\sum\limits_{n=1}^{52} ext{bit}_nrac{1}{2^n}$

Mantisse")

Beim Addieren zweier Zahlen unterschiedlicher Größenordnung geht Information über die kleinere Zahl verloren. Summationen sollten also immer so

 $H_n pprox \log(n) + \gamma + rac{1}{2n} - rac{1}{12n^2} + rac{1}{120n^4}$

%s %s", x, sgn, exponent, mantissa))

bits = bitstring(x) sgn = bits[1]exponent = bits[2:12]

Vorz. Exponent

Mantisse

Mantisse

Mantisse

Die folgende Funktion stellt eine gegebnene Fließkommazahl entsprechend dar.

function maschinendarstellung(x::Float64)

mantissa = bits[13:64]

println(@sprintf("%.15e

println("Dezimal

return nothing

1.110223024625157e-16

maschinendarstellung(-1.0)

maschinendarstellung(1.0+2eps)

-1.000000000000000e+00

1.00000000000000e+00

Out[14]: maschinendarstellung (generic function with 1 method) Schauen wir uns also an wie $1 + \epsilon = 1$ zustande kommt: In [15]: maschinendarstellung(eps)

Exponent

Vorz.

Vorz. Exponent

Vorz. Exponent

durchgeführt werden, dass nur ähnlich große Zahlen miteinander addiert werden. Zurück zur Harmonischen Reihe. Welcher Summationsreihenfolge können wir also trauen? Für großes n gilt

return $log(n)+Base.MathConstants.eulergamma+1.0/(2n)-1.0/(12n^2)+1.0/(120n^4)$

mit der Euler-Gamma Konstante γ . Wir können also unsere beiden Ergebnisse damit vergleichen:

plot(n_werte, abs.(Hn_fwd.-Hn_approx), label="forward", xaxis=:log) plot!(n werte, abs.(Hn bwd.-Hn approx), label="backward", xaxis=:log)

forward backward

0.03 Differenz 0.02 0.01 0.00 10³ 10^{6} n

Out[22]: false

1.32 == (1.2+0.12)

Vergleichen von Fließkommazahlen

der Vergleich innerhalb der numerischen Genauigkeit Sinn.

Prüfen wir zum Beispiel naiv, ob 1.32 gleich 1.2+0.12 ist:

Out[23]: 1.32

Wegen der endlichen Maschinenpräzision müssen wir beim Vergleichen von Fließkommazahlen etwas vorsichtig sein. Dabei ergibt nämlich üblicherweise nur

1.3199999999999998 Zum Verlgeich von Fließkommazahlen innerhalb der numerischen Genauigkeit gibt es die isapprox() Funktion:

In [25]: isapprox(1.32, 1.2+0.12) Out[25]: true