

Scientific Calculations - List 1

Mateusz Pełechaty

23 October 2022

1 Exploration of Arithmetic

1.1 Machine Epsilon

Machine epsilon (*macheps*) is the smallest number x such that $1 + x > 1$ and $rd(1 + x) = 1 + x$.

Find machine epsilon with Julia and compare results with the function `eps` and with values from `float.h`

Solution can be found in

```
./zad1/find_macheps.jl  
./zad1/epsilons.c
```

Method and results

It was calculated by setting $macheps := 1$ and then if $1 + macheps > 1$, then $macheps$ is divided by 2.

Exercise 1.1	Macheps	Eps	Float.h
Float16	0.000977	0.000977	NULL
Float32	1.1920929e-7	1.1920929e-7	1.192093e-07
Float64	2.220446049250313e-16	2.220446049250313e-16	2.220446e-16

Conclusion

Calculating `macheps` by me, `eps(type)` function and `Float.h` constants provide the same values

1.2 Eta

Eta (η) is the smallest number such that $\eta > 0$.

Find η and compare it with `nextfloat(0.0)` and MIN_{sub}

Tests should be made for **Float16**, **Float32**, **Float64**

Solution can be found in

`./zad1/find_eta.jl`

Method and results

It was calculated by setting $\eta := 1$ and then dividing by 2 until $\frac{\eta}{2} > 0$ MIN_{sub} values are taken from *W. Kahan's* book

Exercise 1.2	η	<code>nextfloat</code>	MIN_{sub}
Float16	6.0e-8	6.0e-8	
Float32	1.0e-45	1.0e-45	1.3e-45
Float64	5.0e-324	5.0e-324	4.9e-324

Conclusion

`nextfloat(0.0)` and my method of calculating η provide the same values.

Values are almost the same as MIN_{sub}

1.3 Questions

Q: What is difference between *macheps* and arithmetic precision (ϵ)?

A: *Macheps* is the smallest number that meets condition: $1 + \text{macheps} > 1$.

We can also say that $\text{macheps} = \beta^{1-t}$. ϵ on the other hand is biggest relative error that can happen due to rounding in arithmetic. So it is the smallest number ϵ , that meets condition $\epsilon \geq \delta = \frac{|rd(x)-x|}{x}$ for some number x . It was calculated in the lecture that $\epsilon = \frac{1}{2}\beta^{1-t}$

It follows from here that $\epsilon = \frac{\text{macheps}}{2}$

Q: What is difference between η and (MIN_{sub}) ?

A: MIN_{sub} is minimal subnormal number. η is defined by minimal number bigger than 0. They should be the same, but there are little differences between them **Q:** What is the difference between return value of *floatmin* and

MIN_{nor}

A: They are the same value as seen in table below. Values of MIN_{nor} are taken from *W. Kahan's* book

Q3	floatmin	MIN_{nor}
Float16	6.104e-5	
Float32	1.1754944e-38	1.2E-38
Float64	2.2250738585072014e-308	2.2E-308

1.4 FloatMax

Calculate maximum possible number for Float16, Float32, Float64. Compare values with the ones returned by function *floatmax* and with data

Solution can be found in

`./zad1/floatmax.jl`

Method and results

It was calculated by $max1 := 4$, $max2 := 2$ and $max3 := 1$. Variables are doubled until $max1 == max2$. It means that they are infinity. Then I am returning $max3$

Exercise 1.4	my max	floatmax	W. Kahan's MAX
Float16	3.277e4	6.55e4	
Float32	1.701e38	3.403e38	3.4 E38
Float64	8.988e307	1.798e308	1.8 E308

Conclusion

We can see that *floatmax* is the same as *W. Kahan's* Max

My method of calculating maximum gets me wrong because doubling number reaches infinity. So $mymax = \frac{1}{2} \cdot floatmax$