

Z1/

Epsilon maszynowe:

$\text{fl}(1.0 + \text{eps}) = 1.0 + \text{eps}$

$\text{eps} = 2^{1-t}$

Biasy

Float16 : 15

Float32 : 127

Float64 : 1023

$\text{FloatMax} = +m_{\max} * 2^{c_{\max}}$

$m_{\max} = m_{\max} * 2^0 = 2^1 - \text{eps}(\text{Float})$

$2^{c_{\max}}$ – wyznaczamy iteracyjnie

eps – wyznaczony wcześniej

$\text{FloatMax} = m_{\max} * 2^{c_{\max}} = (2^1 - \text{eps}) * 2^{c_{\max}} = (2^1 - 2^{1-t}) * 2^{c_{\max}}$

EPS for Float64 = 2.220446049250313e-16 ~

EPS for Float32 = 1.1920929e-7 ~

EPS for Float16 = 0.000977 ~

F64:Float64

Macheps64: 2.220446049250313e-16

1.0000000000000002

F32:Float32

Macheps32: 1.1920929e-7

1.0000001

F16:Float16

Macheps16: 0.000977

1.001

NextFloat64: 5.0e-324 ~

NextFloat32: 1.0e-45 ~

NextFloat16: 6.0e-8 ~

E64:Float64

eta: 5.0e-324

E32:Float32

eta: 1.0e-45

E16:Float16

eta: 6.0e-8

MaxFloat64: 1.7976931348623157e308 ~

MaxFloat32: 3.4028235e38 ~

MaxFloat16: 6.55e4 ~

Max 64:Float64

Max: 1.7976931348623157e308

Max 32:Float32

Max: 3.4028235e38

Max 16:Float16

Max: 6.55e4

Float.h:

=== FLOAT ===

FLT_MAX = 3.4028234664e+38

FLT_MIN = 1.1754943508e-38

FLT_EPSILON = 1.1920928955e-07

=== DOUBLE ===

DBL_MAX = 1.79769313486231570815e+308

DBL_MIN = 2.22507385850720138309e-308

DBL_EPSILON = 2.22044604925031308085e-16

=== LONG DOUBLE ===

LDBL_MAX = 1.189731495357231765021263853031e+4932

LDBL_MIN = 3.362103143112093506262677817322e-4932

LDBL_EPSILON = 1.084202172485504434007452800870e-19

Jaki związek ma liczba macheps z precyzją arytmetyki?

Precyzja arytmetyki $\epsilon = 0.5\beta^{1-t}$. Wynika to z tego, że błąd musi spełniać założenia $|\delta| \leq \epsilon$ i $\text{fl}(1.0 + \delta) = 1.0$. Oznacza to, że nie może zostać zaokrąglona w górę, ponieważ wtedy będzie już następną liczbą.

Jaki związek ma liczba eta z liczbą MINsub (zob. wykład lub raport [1])?

$$\text{Min}_{\text{sub}} = m_{\text{min}} * \beta^{\text{cmin}} = \beta^{1-t} * \beta^{\text{cmin}} = \beta^{\text{cmin} + (1-t)}$$

$$\text{MIN}_{\text{nor}} = 1 * \beta^{\text{cmin}}$$

Liczba eta = Minsub

Dlaczego Minsub \neq Macheps

ϵ oznacza precyzję arytmetyki, więc jest najmniejszą liczbą w zakresie cyfr mantysy t. (Tłumaczenie dla mnie żeby zrozumiał)

Co zwracają funkcje floatmin(Float32) i floatmin(Float64) i jaki jest związek zwracanych wartości z liczbą MINnor

Floatmin jest MINnor dla danych typów

biasy- https://en.wikipedia.org/wiki/Exponent_bias

Z2/

eps64: 2.220446049250313e-16 ~

-2.220446049250313e-16

eps32: 1.1920929e-7 ~

1.1920929e-7

eps16: 0.000977 ~

-0.000977

Czasami minus wynika z zaokrągleń

Z3/

$$\delta = 2^{-52}$$

dla [0.5,1] δ jest dwukrotnie za duże

dla [1,2] δ jest idealne

dla [2,4] δ jest dwukrotnie za małe

Z4/

Z5/

Wyniki:

pa: 1.0251881368296672e-10

pb: -1.5643308870494366e-10

pc: 5.653547379013446e6

pd: 5.653547379013446e6

pa32: -0.4999443

pb32: -0.4543457

pc32: 5.653547e6

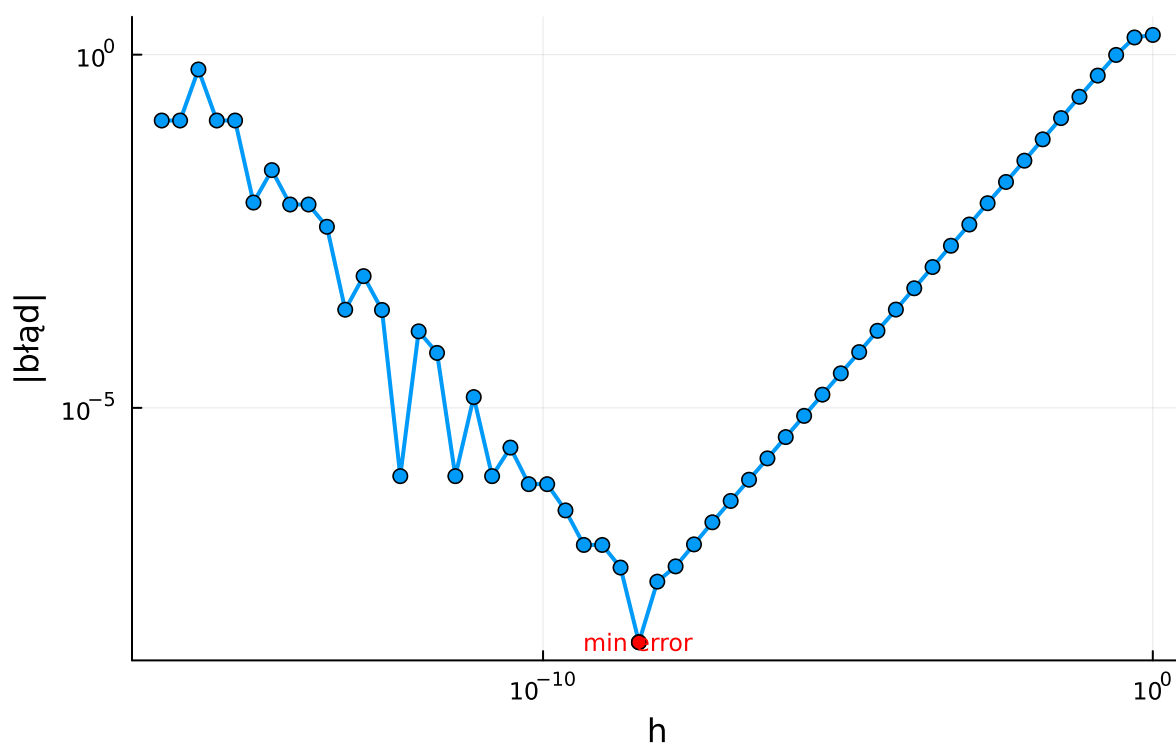
pd32: 5.6535475e6

Z6/

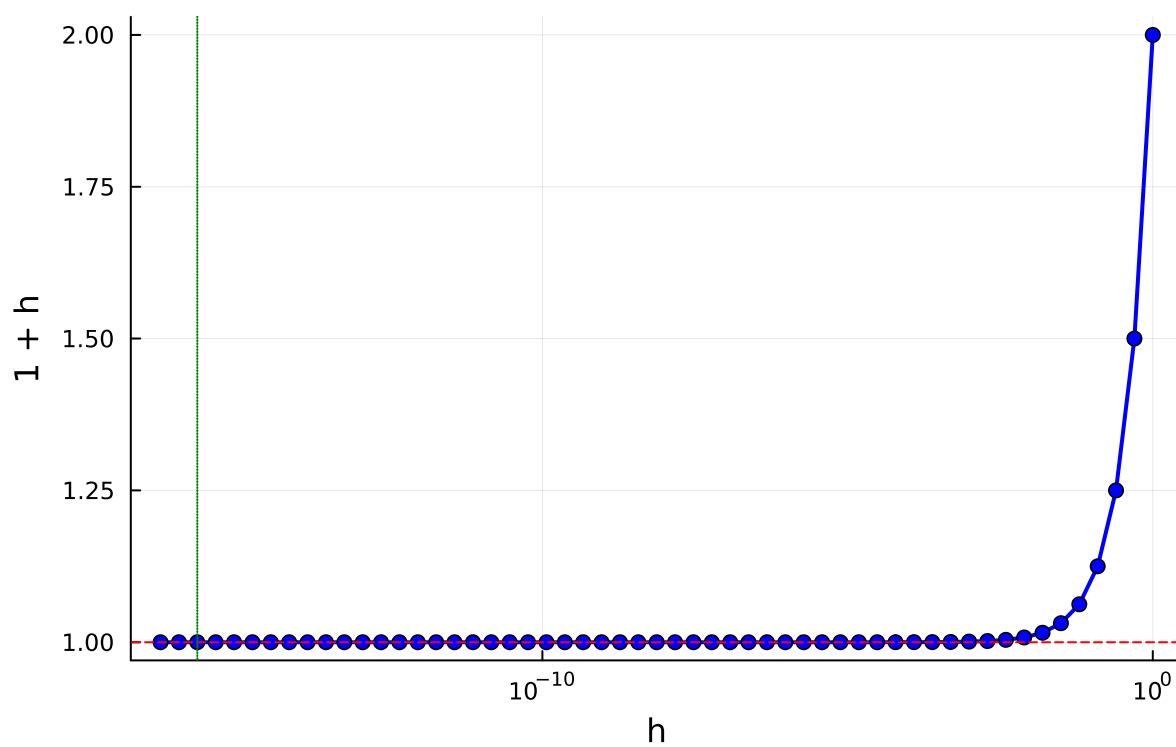
g jest bardziej wiarygodne pewnie przez odejmowanie przy którym musi następować redukcja cyfr znaczących

Z7/

Błąd przybliżenia pochodnej dla $f(x) = \sin(x) + \cos(3x)$



Zachowanie wartości $1 + h$ w arytmetyce Float64



Kod do wykresów wygenerował ChatGPT