

Video 1: Rounding errors

A number system can be represented as $x = \pm 1.b_1b_2b_3b_4 \times 2^m$
for $m \in [-6,6]$ and $b_i \in \{0,1\}$.

Let's say you want to represent the decimal number 19.625 using the binary number system above. Can you represent this number exactly?

$$(19.625)_{10} = (10011.101)_2 = (1.0011101)_2 \times 2^4$$

$$1.0011 \times 2^4 = 19$$

$$1.0100 \times 2^4 = 20$$

Machine floating point number

- Not all real numbers can be exactly represented as a machine floating-point number.

- Consider a real number in the normalized floating point form:

$$x = \pm 1.b_1b_2b_3 \dots b_n \dots \times 2^m$$

$b_{n+1} b_{n+2} \dots$

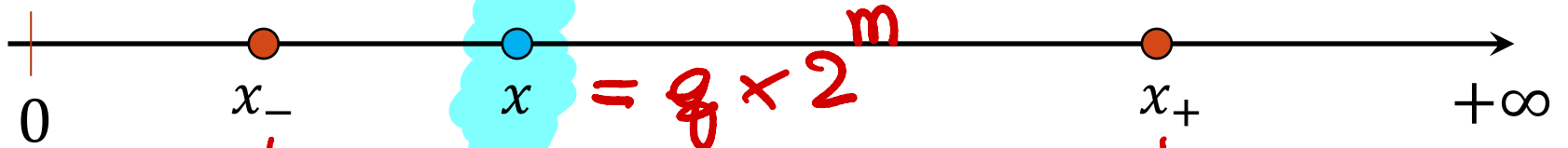
- The real number x will be approximated by either x_- or x_+ , the nearest two machine floating point numbers.



$$x_- = 1.b_1b_2b_3 \dots b_n \times 2^m$$

$$x_+ = x_- + 0.0000 \dots 0 \text{ (1) } \times 2^m$$

$$2^{-n} \times 2^m = \epsilon_m \times 2^m$$



$$x_+ = x_- + \epsilon_m \times 2^m$$

$$|(x_+) - (x_-)| = \epsilon_m \times 2^m$$

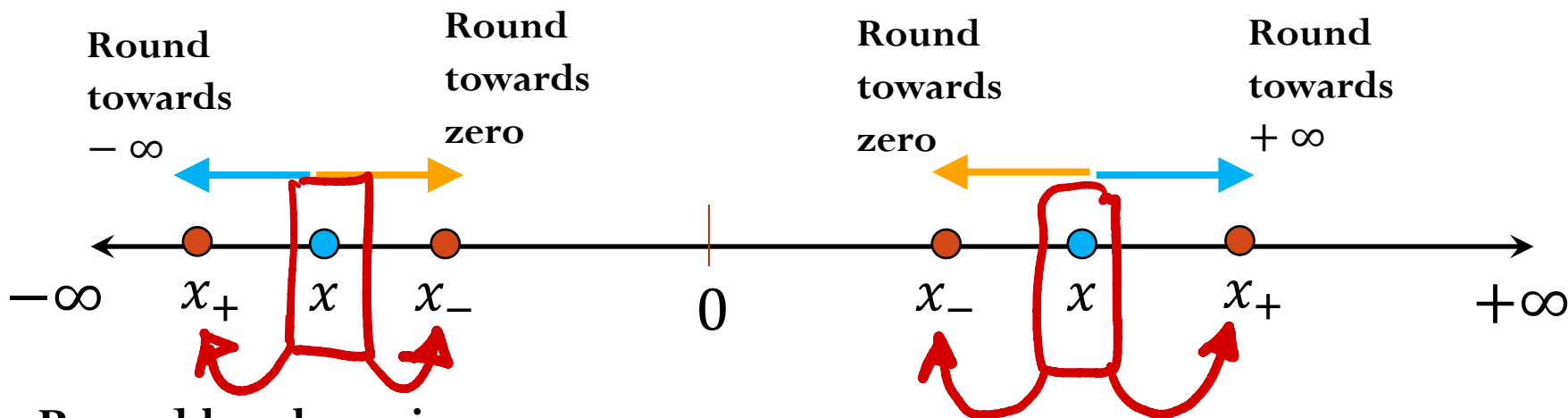
gap = $\epsilon_m \times 2^m$

larger # \longrightarrow larger gap

Rounding

$$x \longrightarrow fl(x) = |fl(x) - x|$$

The process of replacing x by a nearby machine number is called rounding, and the error involved is called **roundoff error**.



Round by chopping:

	x is positive number	x is negative number
Round up (ceil)	$fl(x) = x_+$ Rounding towards $+\infty$	$fl(x) = x_-$ Rounding towards zero
Round down (floor)	$fl(x) = x_-$ Rounding towards zero	$fl(x) = x_+$ Rounding towards $-\infty$

Round to nearest: either round up or round down, whichever is closer

Rounding (roundoff) errors

Consider rounding by chopping:

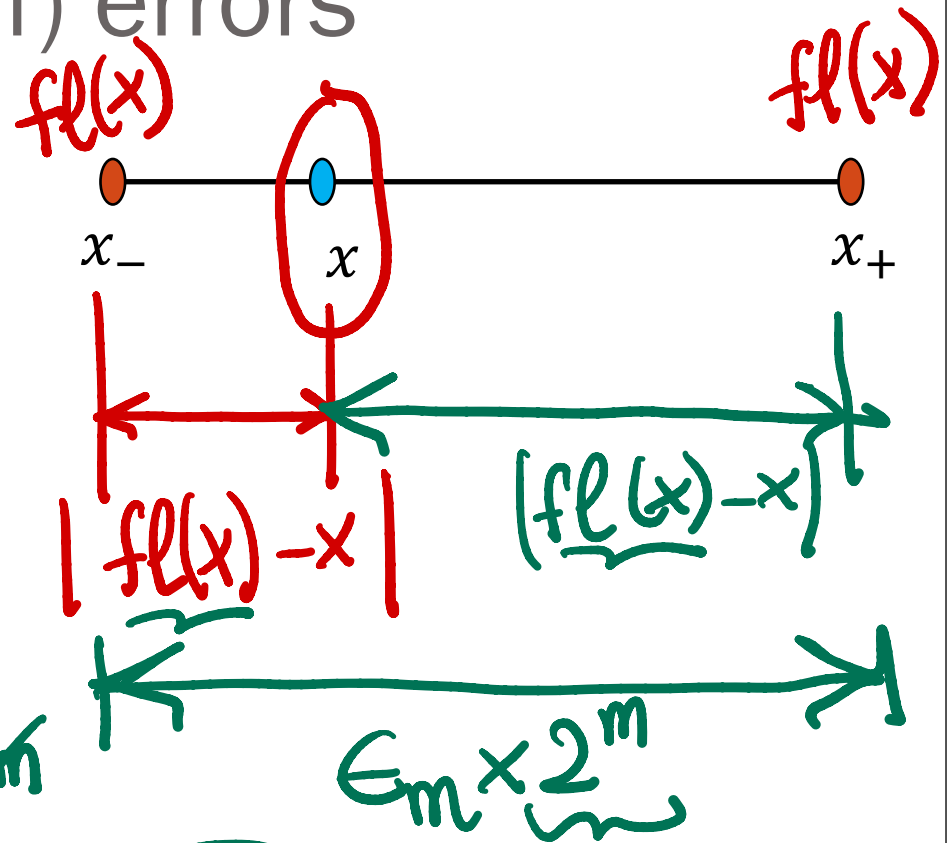
- Absolute error:

$$|fl(x) - x| \leq \epsilon_m \times 2^m$$

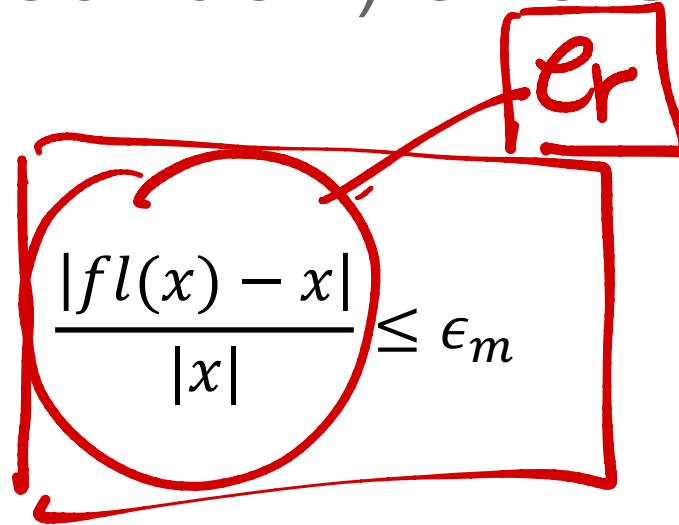
- Relative error:

$$\frac{|fl(x) - x|}{|x|} \leq \frac{\epsilon_m \times 2^m}{|1.b_1b_2 \dots b_n \times 2^m|} = \frac{\epsilon_m}{\underbrace{|1.b_1b_2 \dots b_n|}_{[1, 2)}}$$

$$\boxed{\frac{|fl(x) - x|}{|x|} \leq \epsilon_m}$$



Rounding (roundoff) errors



A hand-drawn diagram in red ink. At the top right, the symbol ϵ_r is enclosed in a small square box. A line from this box points to the fraction in the equation below. The equation is $\frac{|fl(x) - x|}{|x|} \leq \epsilon_m$. The fraction is circled in red, and the entire equation is enclosed in a larger red rectangular box.

$$\frac{|fl(x) - x|}{|x|} \leq \epsilon_m$$

The relative error due to rounding (the process of representing a real number as a machine number) **is always bounded by machine epsilon.**

IEEE Single Precision

$$\frac{|fl(x) - x|}{|x|} \leq 2^{-23} \approx \underline{\underline{1.2 \times 10^{-7}}}$$

$$e_r \leq 1.2 \times 10^{-7} \leq 5 \times 10^{-7}$$

$$e_r \leq 5 \times 10^{(7)}$$

$fl(x) \rightarrow$

7

IEEE Double Precision

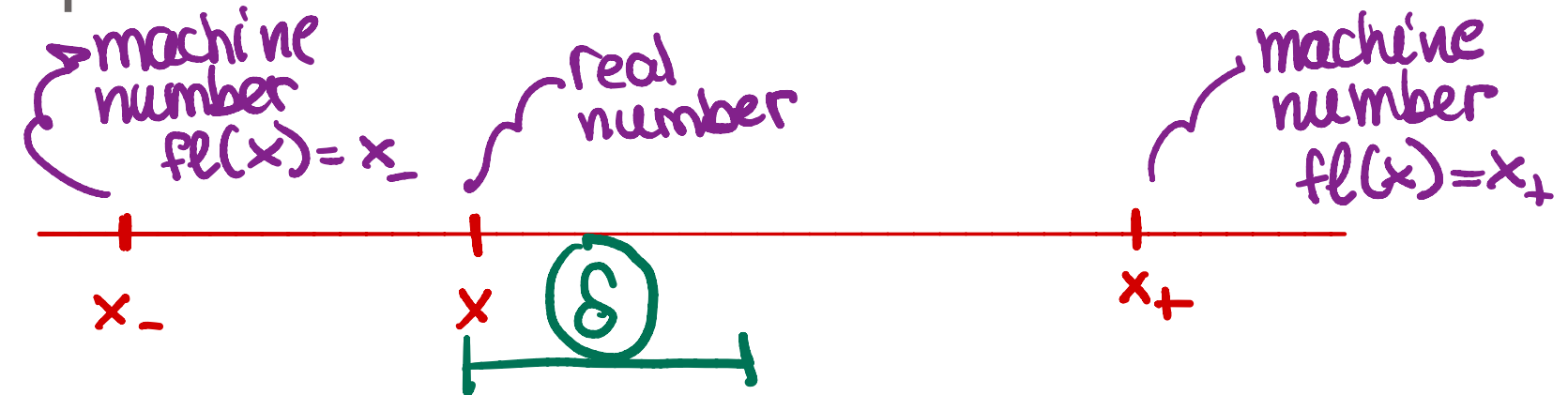
$$\frac{|fl(x) - x|}{|x|} \leq 2^{-52} \approx 2.2 \times 10^{-16}$$

$$e_r \leq 2.2 \times 10^{-16}$$

$$e_r \leq 5 \times 10^{-16}$$

16 decimal

Gap between two machine numbers



$$\hat{x} = fl(x)$$

$$fl(x + \delta) = fl(x) \quad \leftarrow$$

gap

$$\delta < \text{gap}$$

Rule of Thumb

Gap between two machine numbers

$$\text{Binary } x = q \times 2^m$$

$$x = 2^8$$

$$\delta \leq \epsilon_m 2^m \Rightarrow \delta \leq 2^{-23} 2^8$$

(single) $\delta \leq 2^{-15}$

$$\text{fl}(x + \delta) = \text{fl}(x)$$

$$\delta < 2^{-15} \rightarrow \text{fl}(x + \delta) = \text{fl}(x)$$

$$\delta > 2^{-15} \rightarrow \text{fl}(x + \delta) \neq \text{fl}(x)$$

$$\text{Decimal: } x = q \times 10^m$$

$$x = 4.5 \times 10^4$$

double

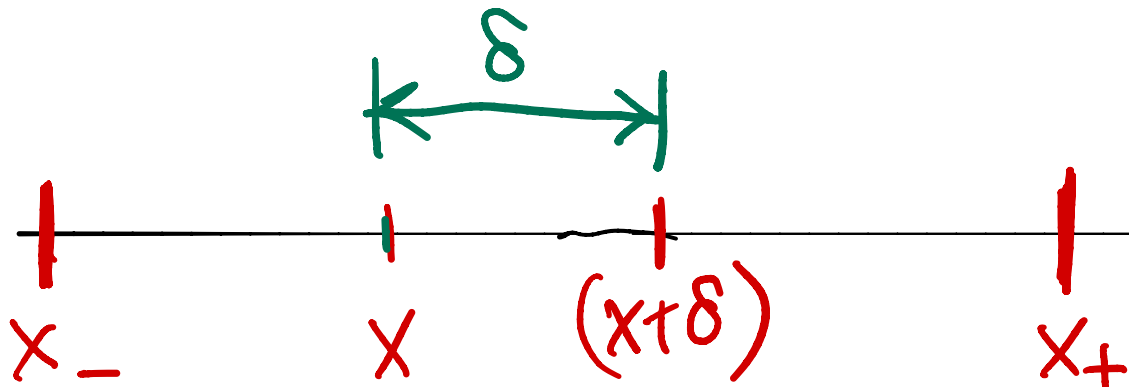
$$\delta \leq \text{gap} \Rightarrow \delta \leq 10^{-16} 10^4$$

$$\delta \leq 10^{-12}$$

$$\delta < 10^{-12} \rightarrow \text{fl}(x + \delta) = \text{fl}(x)$$

$$\delta > 10^{-12} \rightarrow \text{fl}(x + \delta) \neq \text{fl}(x)$$

Gap between two machine numbers



$$fl(x) = \hat{x} = \begin{cases} x_+ \\ x_- \end{cases}$$

$$fl(x+\delta) = \hat{x} = fl(x)$$

What is the smallest δ such that

$$fl(x+\delta) = fl(x) \Rightarrow \delta < gap!$$

$$\delta = gap$$

$$\delta = G_m \times 2^m$$

$$x + \delta$$

$$= q \times 2^m + G_m \times 2^m$$

$$= \underbrace{(q + G_m)}_{\neq q} \times 2^m$$

In practice (Rule of Thumb)
~~Show Python notebook demos~~

Binary base $x = q \times 2^m$

$$fl(x+\delta) = fl(x)$$

$$\delta < \epsilon_m 2^m$$

Example

$$x = 2^8$$

$$\delta < 2^{-23} 2^8 = 2^{-15} \quad \boxed{\delta < 2^{-15}}$$

$$\text{if } \delta < 2^{-15} \Rightarrow fl(x+\delta) = fl(x)$$

$$\text{otherwise } fl(x+\delta) \neq fl(x)$$

Decimal base
 $x = q \times 10^m$

Example $x = 4.5 \times 10^4$

Double Precision

$$\delta < 10^{-16} \times 10^4$$

$$\boxed{\delta < 10^{-12}}$$

Video 2: Arithmetic with machine numbers

Mathematical properties of FP operations

Not necessarily associative:

For some x, y, z the result below is possible:

$$(x + y) + z \neq x + (y + z)$$

```
In [5]: (np.pi+1e100)-1e100
```

```
Out[5]: 0.0
```

```
In [6]: (np.pi)+(1e100-1e100)
```

```
Out[6]: 3.141592653589793
```

```
In [7]: b = 1e80
a = 1e2
print(a + (b - b) )
print((a + b) - b )
```

```
100.0
```

```
0.0
```

Not necessarily distributive:

For some x, y, z the result below is possible:

$$z(x + y) \neq zx + zy$$

```
In [3]: print(100*(0.1 + 0.2))
print(100*0.1 + 100*0.2)
```

```
30.0000000000000004
```

```
30.0
```

```
In [4]: 100*(0.1 + 0.2) == 100*0.1 + 100*0.2
```

```
Out[4]: False
```

Not necessarily cumulative:

Repeatedly adding a very small number to a large number may do nothing

Floating point arithmetic (basic idea)

$$x = (-1)^s 1.f \times 2^m$$

- First compute the exact result
- Then round the result to make it fit into the desired precision
- $x + y = fl(x + y)$
- $x \times y = fl(x \times y)$

Floating point arithmetic

Consider a number system such that $x = \pm 1.b_1b_2b_3 \times 2^m$
for $m \in [-4,4]$ and $b_i \in \{0,1\}$.

$n=3$
 $p=4$

Rough algorithm for addition and subtraction:

1. Bring both numbers onto a common exponent
2. Do “grade-school” operation
3. Round result

- **Example 1: No rounding needed**

$$a = (1.101)_2 \times 2^1$$

$$b = (1.001)_2 \times 2^1$$

$$c = a + b = (10.110)_2 \times 2^1 = (1.011)_2 \times 2^2$$

$$\begin{array}{r} (1.101)_2 \times 2^1 \\ \oplus (1.001)_2 \times 2^1 \\ \hline 10.110 \times 2^1 \end{array}$$

Floating point arithmetic

Consider a number system such that $x = \pm 1.\underbrace{b_1b_2b_3}_{\text{fraction}} \times 2^m$
for $m \in [-4,4]$ and $b_i \in \{0,1\}$.

- **Example 2: Require rounding**

$$\begin{cases} a = (1.101)_2 \times 2^0 \\ b = (1.000)_2 \times 2^0 \end{cases}$$

$$c = a + b = (10.101)_2 \times 2^0$$

$$\longrightarrow 1.0101 \times 2^1$$
$$\text{fl}(a+b) = \underbrace{1.010 \times 2^1}$$

- **Example 3:**

$$\begin{cases} a = (1.100)_2 \times 2^1 \\ b = (1.100)_2 \times 2^{-1} \end{cases}$$

$$c = a + b = \underbrace{(1.100)_2 \times 2^1} + \underbrace{(0.011)_2 \times 2^1} = \underbrace{(1.111)_2 \times 2^1}$$

Floating point arithmetic

Consider a number system such that $x = \pm 1.b_1b_2b_3b_4 \times 2^m$
for $m \in [-4,4]$ and $b_i \in \{0,1\}$.

- **Example 4:**

$$\left\{ \begin{array}{l} a = (1.1011)_2 \times 2^1 \\ b = (1.1010)_2 \times 2^1 \end{array} \right\}$$

$$c = a - b = (0.0001)_2 \times 2^1$$

$$1. \boxed{?} \times 2^1$$

$$fl(a-b) = 1. \underbrace{0000}_1 \times 2^1$$

$n=4 \rightarrow p=5$

$$\begin{array}{r} 1.1011 \times 2^1 \\ - 1.1010 \times 2^1 \\ \hline 0.0001 \times 2^1 \end{array}$$

not signi bits

Floating point arithmetic

Consider a number system such that $x = \pm 1.b_1b_2b_3b_4 \times 2^m$
for $m \in [-4,4]$ and $b_i \in \{0,1\}$.

- **Example 4:**

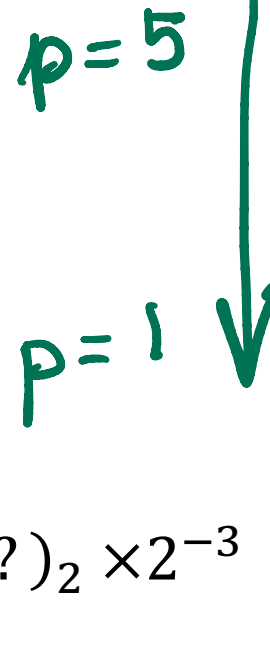
$$a = (1.1011)_2 \times 2^1$$

$$b = (1.1010)_2 \times 2^1$$

$$c = a - b = (0.0001)_2 \times 2^1$$

Or after normalization: $c = (1.????)_2 \times 2^{-3}$

- There is not data to indicate what the missing digits should be.
- Machine fills them with its best guess, which is often not good (usually what is called spurious zeros).
- Number of significant digits in the result is reduced.
- This phenomenon is called **Catastrophic Cancellation**.



Loss of significance

Assume a and b are real numbers with $a \gg b$. For example

$$a = 1.a_1a_2a_3a_4a_5a_6 \dots a_n \dots \times 2^0$$

$$b = 1.b_1b_2b_3b_4b_5b_6 \dots b_n \dots \times 2^{-8}$$

In Single Precision, compute $(a + b)$ $n=23$

$$\begin{aligned} fl(a) &= 1.a_1a_2a_3a_4a_5a_6a_7a_8a_9 \dots a_{22}a_{23} \times 2^0 \\ &\quad 0.00000001b_1b_2 \dots b_{14}b_{15} \times 2^0 \end{aligned}$$

$$fl(a+b) \Rightarrow 15 \text{ bits of } \underline{\underline{b}}$$

Cancellation

Assume a and b are real numbers with $a \approx b$.

$$a = 1.a_1a_2a_3a_4a_5a_6 \dots a_n \dots \times 2^m$$

$$b = 1.b_1b_2b_3b_4b_5b_6 \dots b_n \dots \times 2^m$$

In single precision (without loss of generality), consider this example:

$$a = 1.a_1a_2a_3a_4a_5a_6 \dots a_{20}a_{21}10a_{24}a_{25}a_{26}a_{27} \dots \times 2^m$$

$$b = 1.a_1a_2a_3a_4a_5a_6 \dots a_{20}a_{21}11b_{24}b_{25}b_{26}b_{27} \dots \times 2^m$$

$$b - a = 0.0000 \dots 000 \boxed{1} \times 2^m$$

$$fl(b-a) = 1. \underbrace{000 \dots 00}_{\text{not sig.}} \times 2^{-23} \times 2^m$$

not sig.

Examples:

1) a and b are real numbers with same order of magnitude ($a \approx b$). They have the following representation in a decimal floating point system with 16 decimal digits of accuracy:

$$fl(a) = 3004.45$$

$$fl(b) = 3004.46$$

How many accurate digits does your answer have when you compute $b - a$?

$$\begin{array}{r} 3004.46 \\ - 3004.45 \\ \hline 0000.01 \end{array}$$

11 digits

11 digits

Loss of Significance

How can we avoid this loss of significance? For example, consider the function $f(x) = \sqrt{x^2 + 1} - 1$

If we want to evaluate the function for values x near zero, there is a potential loss of significance in the subtraction.

Assume you are performing this computation using a machine with 5 decimal accurate digits. Compute $f(10^{-3})$

$$f(10^{-3}) = \sqrt{10^{-6} + 1} - 1$$

$\underbrace{\hspace{10em}}_{1 - 1}$

$$= \cancel{\phi}$$

$$\begin{array}{r} 1.000000 \\ + 0.000001 \\ \hline 1.000001 \end{array}$$

Loss of Significance

$$(a-b)(a+b) = a^2 - b^2$$

Re-write the function $f(x) = \sqrt{x^2 + 1} - 1$ to avoid subtraction of two numbers with similar order of magnitude

$$f(x) = (\sqrt{x^2 + 1} - 1) \left(\frac{\sqrt{x^2 + 1} + 1}{\sqrt{x^2 + 1} + 1} \right) = \frac{(\sqrt{x^2 + 1})^2 - (1)^2}{\sqrt{x^2 + 1} + 1}$$

$$= \frac{x^2 + 1 - 1}{\sqrt{x^2 + 1} + 1}$$

$$f(x) = \frac{x^2}{\sqrt{x^2 + 1} + 1}$$

$$f(10^{-3}) = \frac{10^{-6}}{1+1} = \frac{10^{-6}}{2} //$$

round-down

Example:

If $x = 0.3721448693$ and $y = 0.3720214371$ what is the relative error in the computation of $(x - y)$ in a computer with five decimal digits of accuracy?

$$fl(x) = 0.37214$$

$$fl(y) = 0.37202$$

error roundoff

$$\left| \frac{fl(x) - x}{x} \right| = 1.308 \times 10^{-5}$$

$$x - y = 0.0001234322$$

$$fl(x - y) = 0.00012$$

$$e_r = \left| \frac{(x - y) - fl(x - y)}{(x - y)} \right| \approx 3 \times 10^{-2}$$