

# Preserving Significant Digits

## 1 A silly, but telling, example

Consider three functions  $f(x)$ ,  $g(x)$ , and  $h(x)$  defined as

$$f(x) = x + 1 \quad (1)$$

$$g(x) = x \quad (2)$$

$$h(x) = f(x) - g(x) \quad (3)$$

for  $0 \leq x < +\infty$ . Using our keen eyes and razor-sharp intellects, we deduce that  $h(x) = 1$  for all  $x$ . Nonetheless, plotting  $h(x)$  over a large range of  $x$  produces some disturbing results; see Figure 2.

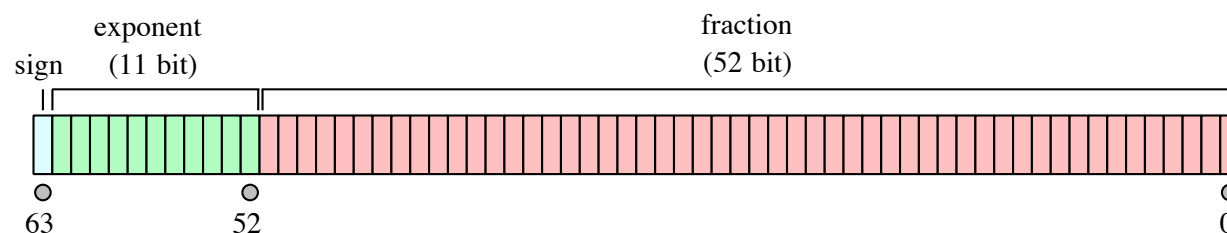
It appears that something unexpected occurs around  $x = 10^{16}$ . Let's zoom in on this part of the graph and see what we can see. Figure 3 reveals the schizophrenic behavior of  $h(x)$ . On the left it is 1, on the right it is 0, and in between it can not seem to make up its mind.

We know why  $h(x) = 1$  on the left – that is what it should be, but two questions remain:

- What is going on around  $x = 10^{16}$ ? It seems to be bouncing between 0 and 2.
- Why is  $h(x) = 0$  for  $x > 10^{17}$ ?

Recall, that a double precision floating point number uses 8 consecutive bytes; that is, 64 bits. The representation has a *sign bit*, 11 bits for the *exponent*, and 52 bits for the *fraction*, as shown in Figure 1.

Figure 1: The IEEE 754 double precision binary floating-point format which is used in most civil engineering applications. This figure was extracted from [en.wikipedia.org/wiki/Double\\_precision](http://en.wikipedia.org/wiki/Double_precision).



Consider the *IEEE 754 double precision floating point number* representation<sup>1</sup> of number  $10^{16}$ :

$$10,000,000,000,000,000 = 0\ 10000110100\ 000111000011011110010011011111000001000000000000000 \quad (4)$$

<sup>1</sup>I used the online calculator at <http://babbage.cs.qc.edu/IEEE-754/Decimal.html> to generate these specific numeric examples.

To reconstruct numbers stored in this format we use

$$\boxed{(-1)^{\text{sign}} \times 2^{\text{exponent}-1023} \times 1.\text{fraction}} \quad (5)$$

Lets parse (4) and verify the representation

sign bit	0
exponent	10000110100 <sub>b</sub> = 1076
fraction	0001110000 1101111001 0011011111 1000001000 0000000000 00 <sub>b</sub>
place	1234567890 1234567890 1234567890 1234567890 1234567890 12

Substituting these pieces of (4) into (5) we have

$$\begin{aligned}
 & (-1)^0 \times 2^{1076-1023} \times 1.000111000011011110010011011111000001000000000000000 \\
 & = 1 \times 2^{53} \times (1 + 2^{-4} + 2^{-5} + 2^{-6} + 2^{-11} + 2^{-12} + \dots + 2^{-37}) \\
 & = 2^{53} + 2^{53-4} + 2^{53-5} + 2^{53-6} + 2^{53-11} + 2^{53-12} + \dots + 2^{53-37} \\
 & = 2^{53} + 2^{49} + 2^{48} + 2^{47} + 2^{42} + 2^{41} + \dots + 2^{16} \\
 & = 10,000,000,000,000,000
 \end{aligned} \quad (6)$$

What is the next larger number, above  $10^{16}$ , that can be represented in the IEEE 754 format?  
To answer this we change the right-most bit of the fraction-part from a 0 to a 1:

sign bit	0
exponent	10000110100 <sub>b</sub> = 1076
fraction	0001110000 1101111001 0011011111 1000001000 0000000000 01 <sub>b</sub>
place	1234567890 1234567890 1234567890 1234567890 1234567890 12

Substituting these pieces into (5) we have

$$\begin{aligned}
 & (-1)^0 \times 2^{1076-1023} \times 1.000111000011011110010011011111000001000000000000001 \\
 & = 1 \times 2^{53} \times (1 + 2^{-4} + 2^{-5} + 2^{-6} + 2^{-11} + 2^{-12} + \dots + 2^{-37} \dots + 2^{-52}) \\
 & = 2^{53} + 2^{53-4} + 2^{53-5} + 2^{53-6} + 2^{53-11} + 2^{53-12} + \dots + 2^{53-37} + \dots + 2^{53-52} \\
 & = 2^{53} + 2^{49} + 2^{48} + 2^{47} + 2^{42} + 2^{41} + \dots + 2^{16} + \dots + 2 \\
 & = 10,000,000,000,000,002
 \end{aligned} \quad (7)$$

There is a gap of 2 in the IEEE 754 format between (10,000,000,000,000,000) and (10,000,000,000,000,002). Thus, the representation for (10,000,000,000,000,000) and (10,000,000,000,000,001) are identical in IEEE 754 format, and

$$\begin{aligned}
 & (10,000,000,000,000,001) - (10,000,000,000,000,000) \\
 & = 0
 \end{aligned} \quad (8)$$

while

$$\begin{aligned}
 & (10,000,000,000,000,002) - (10,000,000,000,000,000) \\
 & = (10,000,000,000,000,002) - (10,000,000,000,000,001) \\
 & = 2
 \end{aligned} \quad (9)$$

## 2 A more challenging example

Consider the function

$$f(x) = \sqrt{x+1} - \sqrt{x} \quad (10)$$

Figure 4 shows what appears to be a smooth function that asymptotically approaches 0, which is what we would expect. Zooming in on the range between  $x = 10^{13}$  and  $10^{17}$  tells a different story; see Figure 5. Are we again suffering from the *slings and arrows of outrageous fortune*? No, it is just the consequences of the truncation error in the IEEE 754 format.

Do we have to simply live with the noise? No. We have to outsmart the noise.

$$f(x) = \sqrt{x+1} - \sqrt{x} \quad (11)$$

$$= (\sqrt{x+1} - \sqrt{x}) \left( \frac{\sqrt{x+1} + \sqrt{x}}{\sqrt{x+1} + \sqrt{x}} \right) \quad (12)$$

$$= \frac{(\sqrt{x+1} - \sqrt{x})(\sqrt{x+1} + \sqrt{x})}{\sqrt{x+1} + \sqrt{x}} \quad (13)$$

$$= \frac{(\sqrt{x+1})^2 - (\sqrt{x})^2}{\sqrt{x+1} + \sqrt{x}} \quad (14)$$

$$= \frac{(x+1) - (x)}{\sqrt{x+1} + \sqrt{x}} \quad (15)$$

$$= \frac{1}{\sqrt{x+1} + \sqrt{x}} \quad (16)$$

Thus, we let

$$g(x) = \frac{1}{\sqrt{x+1} + \sqrt{x}} \quad (17)$$

replace  $f(x)$ , and we slay the dragon.

## 3 The “take home” message

Why are  $h(x)$  in (3), and  $f(x)$  in (10), so badly behaved for large values of  $x$ ? Because both are attempting to evaluate a small difference between very large values of almost identical magnitude.

Be afraid of subtractions between very large values of similar magnitude if you need to maintain significant digits.

On the bright side, in almost all cases, once you see the problem, a mathematical solution is possible.

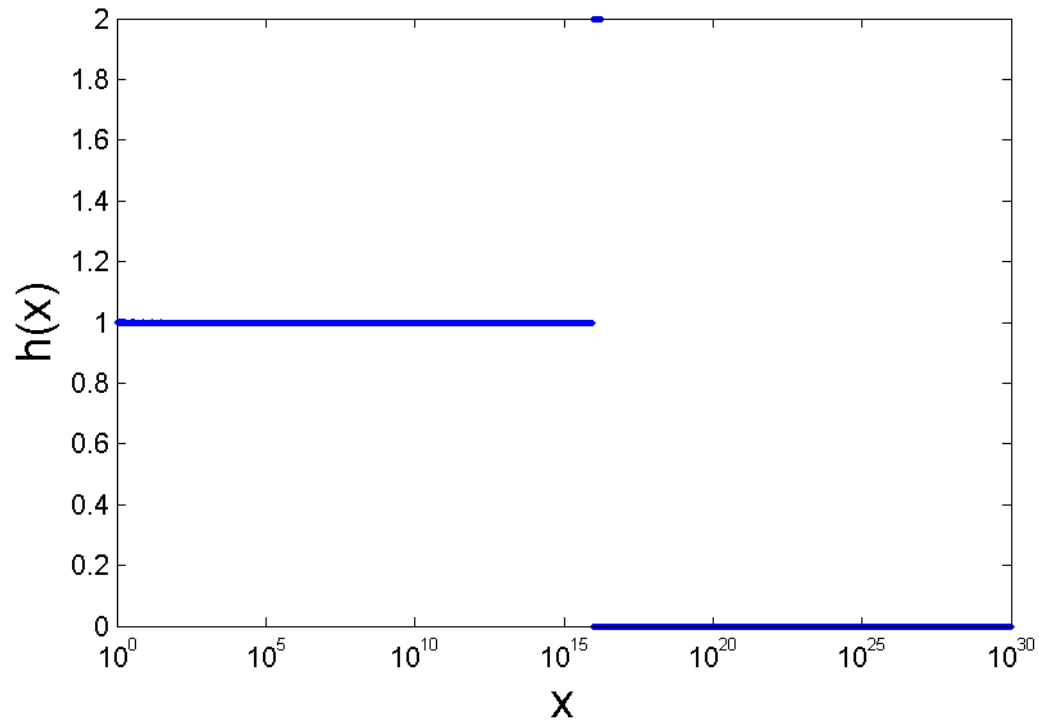


Figure 2: Plotting  $h(x) = (x + 1) - (x)$  over a large range of  $x$  produces disturbing results.

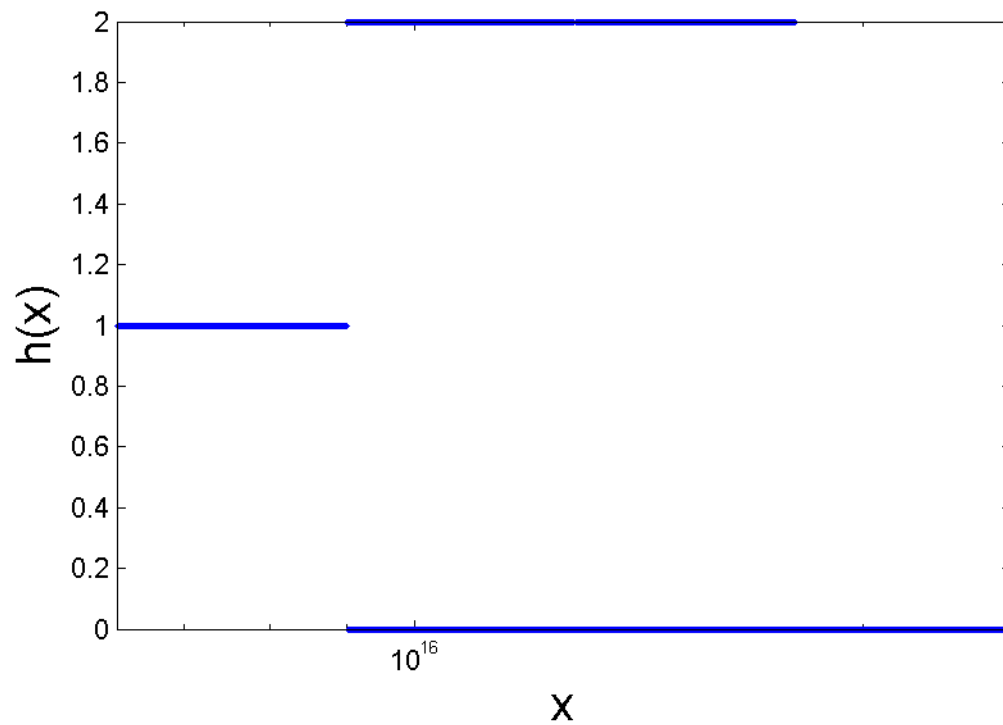


Figure 3: Zooming in around  $x = 10^{16}$  reveals the schizophrenic behavior of  $h(x)$ .

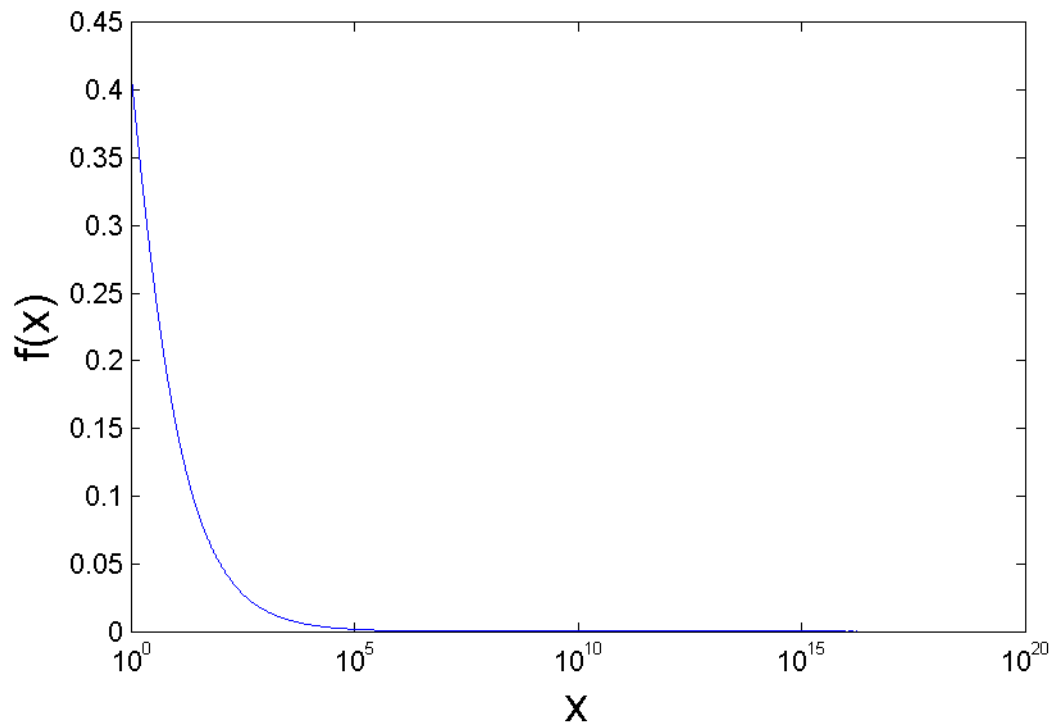


Figure 4: Plotting  $f(x) = \sqrt{x+1} - \sqrt{x}$  over a large range of  $x$ .

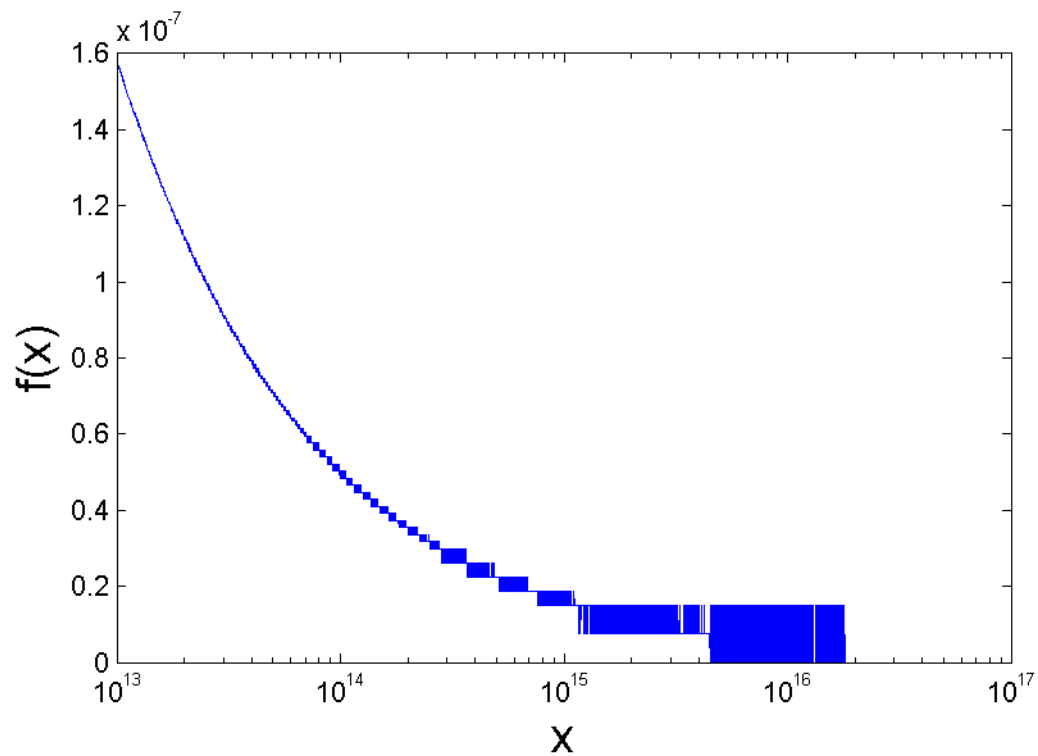


Figure 5: Zooming in between  $x = 10^{13}$  and  $10^{17}$  reveals unexpected misbehavior for  $f(x) = \sqrt{x+1} - \sqrt{x}$ .

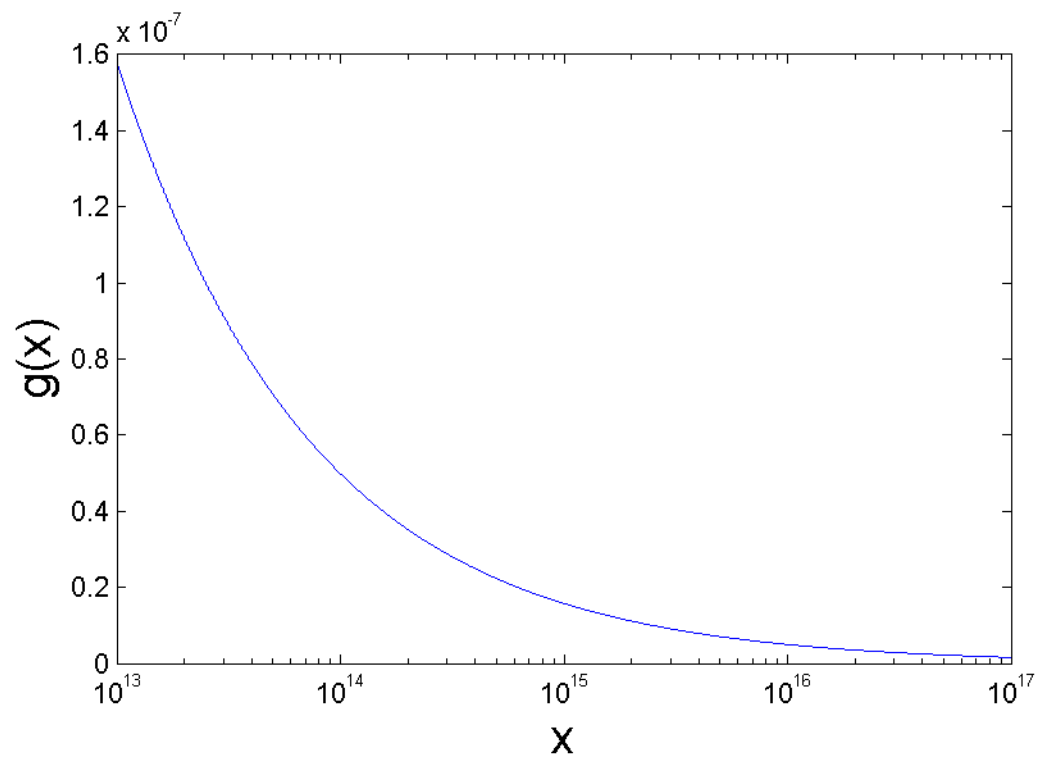


Figure 6: Plotting  $g(x) = \frac{1}{\sqrt{x+1} + \sqrt{x}}$  over the previously troublesome range of  $x$ .