

Math 340 – Fall 2014, Victor Matveev

Binary system, round-off errors, loss of significance, and double precision accuracy.

1. Bits and the binary number system

A “bit” is one digit in a binary representation of a number, and can assume only one of two possible values: **0** or **1**. This simplicity is precisely the reason why the binary number system is more “computer-friendly” than the usual decimal system: in order to store a binary number, it is sufficient to “string together” primitive electric circuits that have only two stable states; one of those states we can call a **0**, and the other state we can call a **1**. In order to store a decimal digit without converting to binary, we would have to use electric circuits with elements that have 10 possible stable states (each state representing a digit from 0 to 9); such a circuit would be extremely complicated to construct, unstable, and completely unnecessary, since binary algebra is easier than decimal algebra.

The binary (base-2) number system is very similar to the usual decimal (base-10) system. It is a **positional** system, which means that the value of each digit comprising the number depends on its position with respect to the fraction point. Namely, in a positional system each digit is multiplied by the system base (10 in the case of decimal, 2 in the case of binary), raised to the power determined by the digit’s position. For instance, you will probably find the following decomposition of number 145.625 to be obvious:

$$\mathbf{145.625} = \mathbf{1} \times 10^2 + \mathbf{4} \times 10^1 + \mathbf{5} \times 10^0 + \mathbf{6} \times 10^{-1} + \mathbf{2} \times 10^{-2} + \mathbf{5} \times 10^{-3}$$

2 1 0 -1 -2 -3

The little number under each digit indicates its position with respect to the fraction point. Now, here is what this number looks like in the binary system:

$$\mathbf{10010001.101}(\text{binary}) = 2^7 + 2^4 + 2^0 + 2^{-1} + 2^{-3} = 128 + 16 + 1 + 0.5 + 0.125 = \mathbf{145.625} \text{ (decimal)}$$

7 6 5 4 3 2 1 0 -1 -2 -3

As this calculation shows, converting from binary to decimal is simple since each binary digit is either **0** or **1**. We don’t have to worry about zeros at all, and it is sufficient to collect powers of two corresponding to each **1**, with the power equal to the position of this **1** with respect to the fraction point.

Now, converting from decimal to binary is only slightly trickier: we have to figure out the highest power of 2 that “fits” within our number, which will give us the position of the highest, most significant “**1**” (the left-most non-zero digit). In the example above, such highest power contained in 145.625 is $128 = 2^7$, so there will be a 1 in position 7. Now, subtract 2^7 from the original number, and find the highest power of 2 that fits in the remainder, 17.625. This will be $16 = 2^4$, so the next **1** is in position 4. Take the remainder, 1.625, and repeat until you are done: $\mathbf{145.625}_{10} = \mathbf{10010001.101}_2$

[Aside: for integer numbers, it is easier to convert from decimal to binary by finding the lowest digit first, instead of the highest one. The lowest digit is equal to the remainder after division by 2]

2. Exponential notation in binary system.

According to the **double precision** standard, **64 bits** are used to represent a real number in computer memory. However, only **52** of those bits are devoted to representing the significant digits, or **mantissa** of the number (also called the **significand**), while another **11** digits are devoted to the “**exponent**”. The remaining one digit determines the sign of the number (it is negative if sign bit equals **1**, and positive otherwise). Now, the “exponent” refers to the standardized exponential scientific notation for real numbers:

$$145.625 = 1.45625 \times 10^2 = 1.45625e+002$$

In this example, the decimal exponent equals two. Re-writing a number in exponential notation is called “normalizing”. In the binary system, each number can be “normalized” in the same way, by carrying the fraction point to the left or to the right. However, for a binary number, each time we move the fraction point by one bit, we pick up a factor of 2, not a factor of 10. For instance, you can easily check that the following is true, by converting both sides of this equation to decimal:

$$\underset{76543210 \ -1-2-3}{10010001.101}(\text{binary}) = 1.0010001101 \times 2^7 \text{ (fraction point moved by 7 bits).}$$

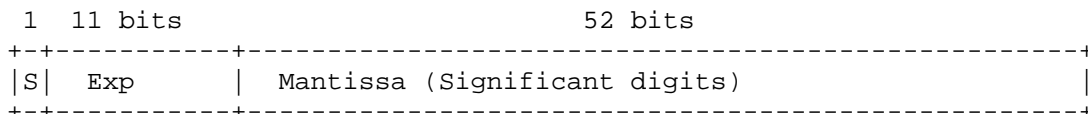
There is a definite purpose of normalizing a number before storing it in the computer: it is easy to see that this way we can store very large or very small numbers, by devoting a certain number of bits (binary digits) to the exponent. Note also that after you normalize a binary number, there will always be a **1** before the fraction point (since **1** is the only non-zero binary digit); therefore, we don’t have to store this **1** – we can simply “pretend” that it’s always there! The only exception is storing zero – we will discuss this below.

3. Double precision: a summary

So sum up, here is a schematic explaining the double precision number representation:

$$\text{Number} = (-1)^S \times 1.F \times 2^{E-1023}$$

Where **S** is the **1-bit** sign, **F** is a **52-bit** long binary mantissa (significand), and **E** is **11** bit-long exponent, which altogether gives **64** bits. In order to represent negative exponent values, 1023 is subtracted from **E** (in principle, one could think of reserving one bit of **E** for the sign, but this subtraction rule allows to define special values of **E** – see below)



Note that the ordering of the bits may be machine-specific: left-to-right, or right-to-left

This shows that double precision representation is very powerful, at least for representing tiny and huge numbers: for comparison, the diameter of our Universe is about 10^{28} - 10^{29} cm, so we rarely encounter numbers bigger or smaller than the above limits!

The real limitation of the double precision standard arises from its accuracy in representing real numbers (the mantissa part, **1.F**). Let's see what the value of the lowest (least significant) 52nd digit in **F** is: $2^{-52}=2.2\times 10^{-16}$. Thus, we have about 16 decimal digits of accuracy to work with. Another way to get the number of significant *decimal* digits in double accuracy is to simply divide **53**, the number of bits, by $\log_2(10)=3.322$, which yields about **16** for the corresponding number of *decimal* places.

5. Round-off errors.

For the sake of simplicity, we will ignore the binary nature of all computer numbers for a moment, and concentrate on the limitations arising from using 16 digits to represent decimal numbers. Consider for instance the following 20-digit integer:

98765432109876543210

Now, this is an example of a number that can not be stored precisely. The large size of the number is not a problem – the problem is the number of significant digits. Since we can only keep 16, we have to “round off” the remaining 4:

98765432109876540000

It is easy to see that the **round-off error** equals 3210 in our example. You can check that Matlab cannot tell the difference between these two numbers: try calculating their difference.

When will this be a problem? If we denote the above number by **N**, you can easily see what the result of the following calculation would be:

$$(N + 1234) - N$$

We would want to get 1234 as our result, but we will naturally get zero instead, since the first four digits are rounded off:

$$98765432109876543210 + 1234 = 98765432109876544444 \rightarrow 98765432109876540000$$

Naturally, the size of the round-off error changes with the size of the number - in other words, if $N=9.8765432109876543210e-109$, then the round-off error will be $3.210e-125$. Using these two examples, you can now understand the following upper bound on the round-off error:

$$|\text{Error of } N| \leq 1\text{e-}16 \times |N|$$

The expression on the right is simply a rough estimate of the 17th digit of N, the one that is thrown out due to round-off.

6. Loss of significance: an example.

Let us examine for instance how the round-off errors may affect a function evaluation, such as $\sin(x + 2\pi n)$. We know that the result has to equal $\sin(x)$. However, if $|x| \leq \text{error}$ of $2\pi n$, there is no way that Matlab (or any other computer program) can recover the value of x from the sum $x + 2\pi n$. In other words, all digits of x will be completely rounded off when x is added to $2\pi n$. This is a classic example of **complete loss of significance of x** . Let's find when this will happen, given $x=1$:

$$x=1 \leq \text{Error}(2\pi n) \leq 1\text{e-}16 \times (2\pi n) \rightarrow n \geq 1\text{e}16 / (2\pi) \approx 1\text{e}15$$

Let us verify this result using Matlab: $\sin(2\pi \times 1\text{e}15 + 1) - \sin(1)$ yields -0.342, which clearly represents complete failure, since this error is comparable to $\sin(1) = 0.841$, while for lower $n = 1\text{e}14$, Matlab gives $\sin(2\pi \times 1\text{e}14 + 1) - \sin(1) = -0.012$, which is still "tolerable".

7. Things to avoid.

As was shown, the round-off errors are usually small (1 part in $1\text{e}16$). Now, apart from the sine calculation example above, when should we in general worry about this error? The answer is simple to formulate: as with any other type of error, **we should worry about the round-off error when it is comparable in absolute value to the result of a calculation that we seek**. In other words, it is not the absolute value of the error which is important, but the *relative* error, which is the ratio of the error to the result value. For instance, here are several such examples where the error is comparable to the result, resulting in a failure when evaluated numerically using Matlab or some other program:

$$(1 + x) - 1 = x \quad (\text{fails when } x \leq 1\text{e-}16)$$

$$(10 - 9)^{14} = \sum_{k=0}^{14} \frac{14!}{k!(14-k)!} 10^k (-9)^{14-k} = 1 \quad (\text{fails})$$

$$\cos(100) = \sum_{n=0}^{\infty} (-1)^n \frac{100^{2n}}{(2n)!} \approx 0.8623 \quad (\text{fails})$$

$$\frac{x^2}{2} \left(\frac{1}{x-1} - \frac{1}{x+1} \right) = \frac{x^2}{x^2 - 1} \quad (\text{fails when } x \geq 1\text{e}16)$$

$$\sqrt{x+1} - \sqrt{x-1} \quad (\text{fails when } x \geq 1\text{e}16)$$

You may not notice at first, but there is something in common between all five examples. Namely, **they all involve precise cancellation between terms to achieve accurate results** (notice that there is a minus sign in each of these formulas). While the error in calculating each of the terms in a given expression may be small compared to the size of that term (1 part in 1e16), it is large when compared to the result, since the result is much smaller than each term taken separately.

Let's take for instance the last example. If our calculation only involves $\sqrt{x+1}$, there is no problem: even though the **1** will get completely rounded off when $x \geq 1e16$, why should we care? The result will be accurate enough; the error relative to the answer will be small (1 part in 1e16, as usual). However, the value of expression $\sqrt{x+1} - \sqrt{x-1}$ is much smaller than each of the two terms taken separately, and when $x \geq 1e16$, yields a zero, which is incorrect. In fact, the exact limiting behavior of this difference is $1/\sqrt{x}$. So the result decreases as x grows, while the round-off error grows; it is easy to see that eventually there will be a problem.

The conclusion is simple: whenever possible, **avoid subtracting numbers that are very close in magnitude!** For instance, the last example may be re-written by multiplying by a conjugate the numerator and the denominator:

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

Although these two formulas are equivalent, the newly derived one is much more accurate when evaluated numerically. For instance, it is easy to see that it gives the correct limiting behavior, $1/\sqrt{x}$, as $x \rightarrow \infty$. Use Matlab to compare these two formulas, for values of x which are greater or equal 1e16.

8. Numbers that are represented exactly, with no round-off errors.

From the earlier discussion of the double precision standard it follows that any binary numbers that can "fit" entirely into the 53 bits allocated for storing its fractional part (which includes the "implied" highest bit) are represented exactly, with no loss of accuracy. Therefore, a number which is given by a sum of powers of two is represented without round-off errors, provided that the difference between the highest and lowest powers is less or equal 52. This includes all integers which are smaller than $2^{53}-1 \approx 9e15$, since such integers can be represented as a binary number with 53 bits or less ($2^{53}-1$ is a binary number with all 53 bits equal to 1: 111111111....1111).

Apart from integers, many fractional numbers are also represented exactly, for instance:

$$\begin{aligned} 0.5 \text{ (decimal)} &= 2^{-1} = 0.1 \text{ (binary)} \\ 0.525 \text{ (decimal)} &= 2^{-1} + 2^{-2} = 0.11 \text{ (binary)} \\ 0.625 \text{ (decimal)} &= 2^{-1} + 2^{-3} = 0.101 \text{ (binary)} \end{aligned}$$

Surprisingly however, some simple decimal numbers which can be exactly represented using only *one decimal* digit, require an *infinite* number of *binary* bits to be represented without error! For instance, you can easily check the validity of the following binary representation of $1/10=0.1$:

$$0.1 \text{ (decimal)} = 0.00011001100110011\dots \text{(binary)} = 2^{-4}+2^{-5}+2^{-8}+2^{-9}+2^{-12}+2^{-13}+\dots$$

The explanation for this fact is very simple: since $0.1 = \frac{1}{2 \times 5}$, it can never be represented as a sum of fractions that have only factors of 2 in the denominator (in other words, 5 is not divisible by 2).

This has serious consequences for calculations involving decimal fractions like 0.1. For example, try the following calculation in Matlab:

$$(0.1 + 0.1 + 0.1 + 0.1 + 0.1 + 0.1 + 0.1 + 0.1) - 0.8$$

To great surprise, it will return $-1.11\text{e-}16$ (which equals -2^{-53} , a round-off error in the lowest bit), not zero! However, if you try

$$(0.5 + 0.5 + 0.5 + 0.5 + 0.5 + 0.5 + 0.5 + 0.5) - 4$$

you will get zero, as expected!