

Numbers on a computer

Chris Paciorek

2023-07-26

Table of contents

1. Basic representations	2
2. Floating point basics	5
Representing real numbers	5
Overflow and underflow	9
Integers or floats?	10
Precision	12
Working with higher precision numbers	16
3. Implications for calculations and comparisons	16
Computer arithmetic is not mathematical arithmetic!	16
Calculating with integers vs. floating points	17
Comparisons	17
Calculations	18
Final note	22

PDF

References:

- Gentle, Computational Statistics, Chapter 2.
- <http://www.lahey.com/float.htm>
- And for more gory detail, see Monahan, Chapter 2.

A quick note that, as we’ve already seen, Python’s version of scientific notation is $X\mathbf{e}Y$, which means $X \cdot 10^Y$.

A second note is that the concepts developed here apply outside of R, but we’ll illustrate the principles of computer numbers using R. R makes use of the *double* and *int* types in C for the underlying representation of R’s numbers in C variables, so what we’ll really be seeing is how such types behave in C on most modern machines. The behavior of real-valued numbers in Python is essentially the same, but Python handles the integer type differently.

Videos (optional):

There are various videos from 2020 in the bCourses Media Gallery that you can use for reference if you want to.

- Video 1. Bits, bytes, and integers
- Video 2. Double precision numbers: intro
- Video 3. Double precision numbers: details
- Video 4. Overflow and integers vs. doubles

1. Basic representations

Everything in computer memory or on disk is stored in terms of bits. A *bit* is essentially a switch that can be either on or off. Thus everything is encoded as numbers in base 2, i.e., 0s and 1s. 8 bits make up a *byte*. For information stored as plain text (ASCII), each byte is used to encode a single character (as previously discussed, actually only 7 of the 8 bits are actually used, hence there are $2^7 = 128$ ASCII characters). One way to represent a byte is to write it in hexadecimal, rather than as 8 0/1 bits. Since there are $2^8 = 256$ possible values in a byte, we can represent it more compactly as 2 base-16 numbers, such as “3e” or “a0” or “ba”. A file format is nothing more than a way of interpreting the bytes in a file.

We’ll create some helper functions to allow us to look at the underlying binary representation.

```
## `Bits` is in the `bitstring` package
def bits(x):
    obj = Bits(float = x, length = 64)
    return(obj.bin)

def bitsi(x, len = 64, hex = False):
    obj = Bits(int = x, length = len)
    if hex:
        return(obj)
    else:
        return(obj.bin)

def dg(x, form = '.20f'):
    print(format(x, form))
```

Note that ‘b’ is encoded as 1 more than ‘a’, and similarly for ‘0’, ‘1’, and ‘2’.

```
Bits(bytes=b'a').bin
```

```
'01100001'
```

```
Bits(bytes=b'b').bin
```

```
'01100010'
```


[illegible]

```
a = 0.3
b = 0.2
dg(a)
```

```
0.2999999999999998890
```

```
dg(b)
```

```
0.20000000000000001110
```

```
dg(a-b)
```

```
0.0999999999999997780
```

```
dg(0.1)
```

```
0.1000000000000000555
```

```
dg(1/3)
```

```
0.3333333333333331483
```

So empirically, it looks like we're accurate up to the 16th decimal place

But actually, the key is the number of digits, not decimal places.

```
dg(1234.1234)
```

```
1234.12339999999994688551
```

```
dg(1234.123412341234)
```

```
1234.12341234123391586763
```

Let's return to our comparison, $0.75 - 0.5 == 0.25$.

```
dg(0.75)
```

```
0.75000000000000000000
```

```
dg(0.50)
```

```
0.50000000000000000000
```

Notice that we can represent the result accurately only up to 16 significant digits. This suggests no need to show more than 16 significant digits and no need to print out any more when writing to

a file (except that if the number is bigger than 10^{16} then we need extra digits to correctly show the magnitude of the number if not using scientific notation). And of course, often we don't need anywhere near that many.

Machine epsilon is the term used for indicating the (relative) accuracy of real numbers and it is defined as the smallest float, x , such that $1 + x \neq 1$:

```
1e-16 + 1.0
1.0

np.array(1e-16) + np.array(1.0)
1.0

1e-15 + 1.0
1.0000000000000001

np.array(1e-15) + np.array(1.0)
1.0000000000000001

2e-16 + 1.0
1.0000000000000002

np.finfo(np.float64).eps
## What about in single precision, e.g. on a GPU?

2.220446049250313e-16

np.finfo(np.float32).eps

1.1920929e-07
```

I'm not sure why numpy reports the the 'resolution' as $1e-15$ when machine epsilon is actually about $2.2e-16$.

Floating point representation

Floating point refers to the decimal point (or *radix* point since we'll be working with base 2 and *decimal* relates to 10). Consider Avogadro's number in terms of scientific notation: $+6.023 \times 10^{23}$. As a baseline for what is about to follow note that we can express a decimal number in the following expansion

$$6.03 = 6 \times 10^0 + 0 \times 10^{-1} + 3 \times 10^{-2}$$

A real number on a computer is stored in what is basically scientific notation:

$$\pm d_0.d_1d_2 \dots d_p \times b^e$$

where b is the base, e is an integer and $d_i \in \{0, \dots, b-1\}$. e is called the *exponent* and $d = d_1 d_2 \dots d_p$ is called the *mantissa*.

Let's consider the choices that the computer pioneers needed to make in using this system to represent numbers on a computer using base 2. First, we need to choose the number of bits to represent e so that we can represent sufficiently large and small numbers. Second we need to choose the number of bits, p , to allocate to $d = d_1 d_2 \dots d_p$, which determines the accuracy of any computer representation of a real.

The great thing about floating points is that we can represent numbers that range from incredibly small to very large while maintaining good precision. The floating point *floats* to adjust to the size of the number. Suppose we had only three digits to use and were in base 10. In floating point notation we can express $0.12 \times 0.12 = 0.0144$ as $(1.20 \times 10^{-1}) \times (1.20 \times 10^{-1}) = 1.44 \times 10^{-2}$, but if we had fixed the decimal point, we'd have $0.120 \times 0.120 = 0.014$ and we'd have lost a digit of accuracy. (Furthermore, we wouldn't be able to represent numbers bigger than 0.99.

More specifically, the actual storage of a number on a computer these days is generally as a double in the form:

$$(-1)^S \times 1.d \times 2^{e-1023} = (-1)^S \times 1.d_1 d_2 \dots d_{52} \times 2^{e-1023}$$

where the computer uses base 2, $b = 2$, (so $d_i \in \{0, 1\}$) because base-2 arithmetic is faster than base-10 arithmetic. The leading 1 normalizes the number; i.e., ensures there is a unique representation for a given computer number. This avoids representing any number in multiple ways, e.g., either $1 = 1.0 \times 2^0 = 0.1 \times 2^1 = 0.01 \times 2^2$. For a double, we have 8 bytes=64 bits. Consider our representation as (S, d, e) where S is the sign. The leading 1 is the *hidden bit* and doesn't need to be stored because it is always present. In general e is represented using 11 bits ($2^{11} = 2048$), and the subtraction takes the place of having a sign bit for the exponent. (Note that in our discussion we'll just think of e in terms of its base 10 representation, although it is of course represented in base 2.) This leaves $p = 52 = 64 - 1 - 11$ bits for d .

In this code I'll force storage as a double using by tacking on a decimal place, .0.

```
bits(2.0**(-1)) # 1/2
```

```
'00111111110000000000000000000000000000000000000000000000000000'
```

```
bits(2.0**0) # 1
```

```
'00111111111000000000000000000000000000000000000000000000000000000000'
```

```
bits(2.0**1) # 2
```

[illegible]

```
bits(2.0**1 + 2.0**0) # 3
```


[illegible]

Let's consider what can be represented exactly:

```
dg(.5)
0.50000000000000000000
```

```
dg(.26)
0.26000000000000000888
```

dg(1/33)

0.030303030303030387

Overflow and underflow

9

infer about the range of possible numbers. With 11 bits for e , we can represent $\pm 2^{10} = \pm 1024$ different exponent values (see `np.finfo(np.float64).maxexp`) (why is `np.finfo(np.float64).minexp` only -1022?). So the largest number we could represent is 2^{1024} . What is this in base 10?

```
1e308

1e+308

1e309

inf

try:
    np.log10(2.0**1024)
except:
    print("Whoops -- that just barely overflowed.")
```

Whoops -- that just barely overflowed.

```
np.log10(2.0**1023)
```

```
307.95368556425274
```

```
np.finfo(np.float64)
```

```
finfo(resolution=1e-15, min=-1.7976931348623157e+308, max=1.7976931348623157e+308, dtype=float64)
```

We could have been smarter about that calculation: $\log_{10} 2^{1024} = \log_2 2^{1024} / \log_2 10 = 1024 / 3.32 \approx 308$. The result is analogous for the smallest number, so we have that floating points can range between 1×10^{-308} and 1×10^{308} . Take a look at `.Machine$double.xmax` and `.Machine.double.xmin`. Producing something larger or smaller in magnitude than these values is called overflow and underflow respectively. When we overflow, R gives back an Inf or -Inf (and in other cases we might get an error message). When we underflow, we get back 0, which in particular can be a problem if we try to divide by the value.

Integers or floats?

Values stored as integers should overflow if they exceed the maximum integer.

Should 2^{65} overflow?

```
np.log2(np.iinfo(np.int64).max)
```

```
63.0
```

```

x = np.int64(2)
# Yikes!
x**64

# Interesting:

0

2**64
18446744073709551616

2**100
1267650600228229401496703205376

dg(2.0**64, '.2f')
18446744073709551616.00

dg(2.0**100, '.2f')
1267650600228229401496703205376.00

```

Doubles won't overflow until much larger values than 4- or 8-byte integers. However we need to think about what integer-valued numbers can and can't be stored exactly in our base 2 representation of floating point numbers. It turns out that integer-valued numbers can be stored exactly as doubles when their absolute value is less than 2^{53} .

Challenge: Why 2^{53} ? Write out what integers can be stored exactly in our base 2 representation of floating point numbers.

You can force storage as integers or doubles in a few ways.

```

x = 3; type(x)
<class 'int'>

x = np.float64(x); type(x)
<class 'numpy.float64'>

x = 3.0; type(x)
<class 'float'>

```

```
x = np.float64(3); type(x)
```

```
<class 'numpy.float64'>
```

Precision

Consider our representation as (S, d, e) where we have $p = 52$ bits for d . Since we have $2^{52} \approx 0.5 \times 10^{16}$, we can represent about that many discrete values, which means we can accurately represent about 16 digits (in base 10). The result is that floats on a computer are actually discrete (we have a finite number of bits), and if we get a number that is in one of the gaps (there are uncountably many reals), it's approximated by the nearest discrete value. The accuracy of our representation is to within $1/2$ of the gap between the two discrete values bracketing the true number. Let's consider the implications for accuracy in working with large and small numbers. By changing e we can change the magnitude of a number. So regardless of whether we have a very large or small number, we have about 16 digits of accuracy, since the absolute spacing depends on what value is represented by the least significant digit (the *ulp*, or *unit in the last place*) in d , i.e., the $p = 52$ nd one, or in terms of base 10, the 16th digit. Let's explore this:

```
# large vs. small numbers
dg(.1234123412341234)
```

```
0.12341234123412339607
```

```
dg(1234.1234123412341234) # not accurate to 16 decimal places
```

```
1234.12341234123414324131
```

```
dg(123412341234.123412341234) # only accurate to 4 places
```

```
123412341234.12341308593750000000
```

```
dg(1234123412341234.123412341234) # no places!
```

```
1234123412341234.00000000000000000000
```

```
dg(12341234123412341234) # fewer than no places!
```

```
12341234123412340736.00000000000000000000
```

We can see the implications of this in the context of calculations:

```
dg(1234567812345678.0 - 1234567812345677.0)
```

```
1.000000000000000000000000
```

```
dg(12345678123456788888.0 - 12345678123456788887.0)
```

```
0.00000000000000000000
```

```
dg(12345678123456780000.0 - 12345678123456770000.0)
```

```
10240.00000000000000000000
```

The spacing of possible computer numbers that have a magnitude of about 1 leads us to another definition of *machine epsilon* (an alternative, but essentially equivalent definition to that given previously in this Unit). Machine epsilon tells us also about the relative spacing of numbers. First let's consider numbers of magnitude one. The difference between $1 = 1.00...00 \times 2^0$ and $1.000...01 \times 2^0$ is $\epsilon = 1 \times 2^{-52} \approx 2.2 \times 10^{-16}$. Machine epsilon gives the *absolute spacing* for numbers near 1 and the *relative spacing* for numbers with a different order of magnitude and therefore a different absolute magnitude of the error in representing a real. The relative spacing at x is

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

since the next largest number from x is given by $(1 + \epsilon)x$.

Suppose $x = 1 \times 10^6$. Then the absolute error in representing a number of this magnitude is $x\epsilon \approx 2 \times 10^{-10}$. (Actually the error would be one-half of the spacing, but that's a minor distinction.) We can see by looking at the numbers in decimal form, where we are accurate to the order 10^{-10} but not 10^{-11} . This is equivalent to our discussion that we have only 16 digits of accuracy.

```
dg(1000000.1)
```

```
1000000.09999999997671693563
```

Let's see what arithmetic we can do exactly with integer-valued numbers stored as doubles and how that relates to the absolute spacing of numbers we've just seen:

```
2.0**52
```

```
4503599627370496.0
```

```
2.0**52+1
```

```
4503599627370497.0
```

```
2.0**53
```

```
9007199254740992.0
```

```
2.0**53+1
```

```
9007199254740992.0
```


for numbers around 0.1. The numbers R reports are spaced in increments of individual bits in the base 2 representation.

```
dg(0.1234567812345678)
0.12345678123456779729

dg(0.12345678123456781)
0.12345678123456781117

dg(0.12345678123456782)
0.12345678123456782505

dg(0.12345678123456783)
0.12345678123456782505

dg(0.12345678123456784)
0.12345678123456783892

bits(0.1234567812345678)
'001111110111111100110101101110100010101110111110011010010000110'

bits(0.12345678123456781)
'001111110111111100110101101110100010101110111110011010010000111'

bits(0.12345678123456782)
'001111110111111100110101101110100010101110111110011010010001000'

bits(0.12345678123456783)
'001111110111111100110101101110100010101110111110011010010001000'

bits(0.12345678123456784)
'001111110111111100110101101110100010101110111110011010010001001'
```

Working with higher precision numbers

As we've seen, Python will automatically work with integers in arbitrary precision. (Note that R does not do this – R uses 4-byte integers, and for many calculations it's best to use R's `numeric` type because integers that aren't really large can be expressed exactly.)

For higher precision floating point numbers you can make use of the `gmpy2` package.

```
import gmpy2
gmpy2.get_context().precision=200
gmpy2.const_pi()

## not sure why this shows ...00004
gmpy2.mpfr(".1234567812345678")
```

3. Implications for calculations and comparisons

Computer arithmetic is not mathematical arithmetic!

As mentioned for integers, computer number arithmetic is not closed, unlike real arithmetic. For example, if we multiply two computer floating points, we can overflow and not get back another computer floating point. One term that is used, which might pop up in an error message (though probably not in R) is that an “exception” is “thrown”.

Another mathematical concept we should consider here is that computer arithmetic does not obey the associative and distributive laws, i.e., $(a + b) + c$ may not equal $a + (b + c)$ on a computer and $a(b + c)$ may not be the same as $ab + ac$. Here's an example with multiplication:

```
val1 = 1/10; val2 = 0.31; val3 = 0.57
res1 = val1*val2*val3
res2 = val3*val2*val1
res1 == res2
```

False

```
dg(res1)
```

0.01766999999999999821

```
dg(res2)
```

0.01767000000000000168

Calculating with integers vs. floating points

It's important to note that operations with integers are fast and exact (but can easily overflow) while operations with floating points are slower and approximate. Because of this slowness, floating point operations (*flops*) dominate calculation intensity and are used as the metric for the amount of work being done - a multiplication (or division) combined with an addition (or subtraction) is one flop. We'll talk a lot about flops in the unit on linear algebra.

Comparisons

As we saw, we should never test `a==b` unless (1) *a* and *b* are represented as integers in R, (2) they are integer-valued but stored as doubles that are small enough that they can be stored exactly) or (3) they are decimal numbers that have been created in the same way (e.g., `0.4-0.3==0.4-0.3` returns `TRUE` but `0.1==0.4-0.3` does not). Similarly we should be careful about testing `a==0`. And be careful of greater than/less than comparisons. For example, be careful of `x[x < 0] = NA` if what you are looking for is values that might be *mathematically* less than zero, rather than whatever is *numerically* less than zero.

```
4 - 3 == 1
```

True

```
4.0 - 3.0 == 1.0
```

True

```
4.1 - 3.1 == 1.0
```

False

One nice approach to checking for approximate equality is to make use of *machine epsilon*. If the relative spacing of two numbers is less than *machine epsilon*, then for our computer approximation, we say they are the same. Here's an implementation that relies on the absolute spacing being $x\epsilon$ (see above).

```
a = 12345678123456781000
b = 12345678123456782000

def approxEqual(a,b):
    if abs(a - b) < np.finfo(np.float64).eps * abs(a + b):
        print("approximately equal")
    else:
        print ("not equal")

approxEqual(a,b)
```

```
a = 1234567812345678
b = 1234567812345677

approxEqual(a,b)
```

Actually, we probably want to use a number slightly larger than machine epsilon to be safe.

Calculations

1. Subtracting large numbers that are nearly equal (or adding negative and positive numbers of the same magnitude). You won't have the precision in the answer that you would like. How many decimal places of accuracy do we have here?

The absolute error in the original numbers here is of the order $\epsilon x = 2.2 \times 10^{-16} \cdot 1 \times 10^{11} \approx 1 \times 10^{-5} = .00001$. While we might think that the result is close to the value 1 and should have error of about machine epsilon, the relevant absolute error is in the original numbers, so we actually only have about five significant digits in our result because we cancel out the other digits.

```
# catastrophic cancellation w/ small numbers
a = .000000000000123412341234
b = .000000000000123412340000

# so we know the right answer is .000000000000000000001234 EXACTLY
```

0.000000000000000000000000123399999315140

It's best to do any subtraction on numbers that are not too large. For example, if we compute the sum of squares in a naive way, we can lose all of the information in the calculation because the information is in digits that are not computed or stored accurately:

```
## No problem here:
x = np.array([-1.0, 0.0, 1.0])
n = len(x)
np.sum(x**2)-n*np.mean(x)**2
```

```
np.sum((x - np.mean(x))**2)
```

2.0

0.0

2.0

19

```

dg(123456781234.2)
123456781234.19999694824218750000

dg(123456781234.2 - 0.1)      # truth: 123456781234.1
123456781234.09999084472656250000

dg(123456781234.2 - 0.01)     # truth: 123456781234.19
123456781234.19000244140625000000

dg(123456781234.2 - 0.001)    # truth: 123456781234.199
123456781234.19898986816406250000

dg(123456781234.2 - 0.0001)   # truth: 123456781234.1999
123456781234.19989013671875000000

dg(123456781234.2 - 0.00001)  # truth: 123456781234.19999
123456781234.19998168945312500000

dg(123456781234.2 - 0.000001) # truth: 123456781234.199999
123456781234.19999694824218750000

123456781234.2 - 0.000001 == 123456781234.2

```

True

The larger number in the calculations above is of magnitude 10^{11} , so the absolute error in representing the larger number is around 1×10^{-5} . Thus in the calculations above we can only expect the answers to be accurate to about 1×10^{-5} . In the last calculation above, the smaller number is smaller than 1×10^{-5} and so doing the subtraction has had no effect. This is analogous to trying to do $1 + 1 \times 10^{-16}$ and seeing that the result is still 1.

A work-around when we are adding numbers of very different magnitudes is to add a set of numbers in increasing order. However, if the numbers are all of similar magnitude, then by the time you add ones later in the summation, the partial sum will be much larger than the new term. A (second) work-around to that problem is to add the numbers in a tree-like fashion, so that each addition involves a summation of numbers of similar size.

Given the limited *range* of computer numbers, be careful when you are:

- Multiplying or dividing many numbers, particularly large or small ones. Never take the product of many large or small numbers as this can cause over- or under-flow. Rather compute on the log scale and only at the end of your computations should you exponentiate. E.g.,

$$\prod_i x_i / \prod_j y_j = \exp(\sum_i \log x_i - \sum_j \log y_j)$$

Let's consider some challenges that illustrate that last concern.

- Challenge: consider multiclass logistic regression, where you have quantities like this:

$$p_j = \text{Prob}(y = j) = \frac{\exp(x\beta_j)}{\sum_{k=1}^K \exp(x\beta_k)} = \frac{\exp(z_j)}{\sum_{k=1}^K \exp(z_k)}$$

for $z_k = x\beta_k$. What will happen if the z values are very large in magnitude (either positive or negative)? How can we reexpress the equation so as to be able to do the calculation? Hint: think about multiplying by $\frac{c}{c}$ for a carefully chosen c .

- Second challenge: The same issue arises in the following calculation. Suppose I want to calculate a predictive density (e.g., in a model comparison in a Bayesian context):

$$\begin{aligned} f(y^*|y, x) &= \int f(y^*|y, x, \theta) \pi(\theta|y, x) d\theta \\ &\approx \frac{1}{m} \sum_{j=1}^m \prod_{i=1}^n f(y_i^*|x, \theta_j) \\ &= \frac{1}{m} \sum_{j=1}^m \exp \sum_{i=1}^n \log f(y_i^*|x, \theta_j) \\ &\equiv \frac{1}{m} \sum_{j=1}^m \exp(v_j) \end{aligned}$$

First, why do I use the log conditional predictive density? Second, let's work with an estimate of the unconditional predictive density on the log scale, $\log f(y^*|y, x) \approx \log \frac{1}{m} \sum_{j=1}^m \exp(v_j)$. Now note that e^{v_j} may be quite small as v_j is the sum of log likelihoods. So what happens if we have terms something like e^{-1000} ? So we can't exponentiate each individual v_j . This is what is known as the "log sum of exponentials" problem (and the solution as the "log-sum-exp trick"). Thoughts?

Numerical issues come up frequently in linear algebra. For example, they come up in working with positive definite and semi-positive-definite matrices, such as covariance matrices. You can easily get negative numerical eigenvalues even if all the eigenvalues are positive or non-negative. Here's an example where we use an squared exponential correlation as a function of time (or distance in 1-d), which is *mathematically* positive definite (i.e., all the eigenvalues are positive) but not numerically positive definite:

```
def custom_operation(a, b):
    return abs(a - b) # Custom pairwise operation
```

```

xs = np.arange(100)
dists = custom_operation(xs[:, np.newaxis], xs)
corMat = np.exp(-(dists/10)**2) # this is a p.d. matrix (mathematically)
scipy.linalg.eigvals(corMat)[80:99] # but not numerically

array([ 1.87147120e-16-1.38569039e-16j, -2.63753999e-16+0.00000000e+00j,
       2.02746697e-16+9.27811726e-17j,  2.02746697e-16-9.27811726e-17j,
      -2.18770059e-16+0.00000000e+00j, -1.81043922e-16+6.92659199e-17j,
      -1.81043922e-16-6.92659199e-17j,  2.16131809e-16+2.53937305e-17j,
       2.16131809e-16-2.53937305e-17j, -1.48286621e-16+5.22080636e-17j,
      -1.48286621e-16-5.22080636e-17j, -1.45202652e-16+0.00000000e+00j,
      -5.96180427e-18+1.51715979e-16j, -5.96180427e-18-1.51715979e-16j,
      -4.10075792e-17+5.46543960e-17j, -4.10075792e-17-5.46543960e-17j,
       1.51314597e-16+0.00000000e+00j,  5.42455332e-17+4.22706275e-17j,
       5.42455332e-17-4.22706275e-17j])

```

Final note

How the computer actually does arithmetic with the floating point representation in base 2 gets pretty complicated, and we won't go into the details. These rules of thumb should be enough for our practical purposes. Monahan and the URL reference have many of the gory details.