

Chapter 1

Fun with Floats

with the participation of:
Nicolas Cellier (*ncellier@ifrance.com*)

Floats are inexact by nature and this can confuse programmers. In this chapter we present an introduction to this problem. The basic message is that Floats are what they are: inexact but fast numbers.

1.1 Never test equality on floats

The first basic principle is to never compare float equality. Let's take a simple case: the addition of two floats may not be equal to the float representing their sum. For example $0.1 + 0.2$ is not equal to 0.3 .

```
(0.1 + 0.2) = 0.3  
returns false
```

Hey, this is unexpected, you did not learn that in school, did you? This behavior is surprising indeed, but it's normal since floats are inexact numbers. What is important to understand is that the way floats are printed is also impacting our understanding. Some approaches prints a simpler representation of reality than others. In early versions of Pharo printing $0.1 + 0.2$ were printing 0.3 , now printing it returns 0.30000000000000004 . This change was guided by the idea that it is better not to lie to the user. Showing the inexactness of float is better than hiding because one day or another we can be deeply bitten by them.

```
(0.2 + 0.1) printString  
returns '0.30000000000000004'
```

```
0.3 printString
returns '0.3'
```

We can see that we are in presence of two different numbers by looking at the hexadecimal values.

```
(0.1+0.2) hex
returns '3FD3333333333334'
0.3 hex
returns '3FD3333333333333'
```

The method `storeString` also conveys that we are in presence of two different numbers.

```
(0.1+0.2) storeString
returns '0.30000000000000004'
0.3 storeString
returns '0.3'
```

About `closeTo`: One way to know if two floats are probably close enough to look like the same number is to use the message `closeTo`:

```
(0.1 + 0.2) closeTo: 0.3
returns true

0.3 closeTo: (0.1 + 0.2)
returns true
```

About Scaled Decimals. There is a solution if you absolutely need exact floating point numbers, use Scaled Decimals. Scaled Decimals are exact numbers so they exhibit the behavior you expected.

```
0.1s2 + 0.2s2 = 0.3s2
returns true
```

1.2 Disecting a Float

To understand what operation is involved in above addition, we must know how Floats are represented internally in the computer: Pharo's Float format is a wide spread standard found on most computers - IEEE 754-1985 double precision on 64 bits (See http://en.wikipedia.org/wiki/IEEE_754-1985 for more details). In this format, a Float is represented in base 2 by this formula:

$$sign \cdot mantissa \cdot 2^{exponent}$$


```
returns 53
```

```
Float precision.  
returns 53
```

You can also retrieve the exact fraction corresponding to the internal representation of the Float:

```
11.125 asTrueFraction.  
returns (89/8)
```

```
(#(0 2 3 6) detectSum: [:i | (2 raisedTo: i) reciprocal]) * (2 raisedTo: 3).  
returns (89/8)
```

Until there we retrieved the exact input we've injected into the Float. Are Float operations exact after all? Hem, no, we only played with fractions having a power of 2 as denominator and a few bits in numerator. If one of these conditions is not met, we won't find any exact Float representation of our numbers. In particular, it is not possible to represent $1/5$ with a finite number of binary digits. Consequently, a decimal fraction like 0.1 cannot be represented exactly with above representation.

```
(1/5) asFloat = (1/5).  
returns false
```

```
(1/10) = 0.1  
returns false
```

Let us see in detail how we could get the fractional bits of $1/10$ *i.e.*, $2r1 / 2r101$. For that, we must pose the division:

1	101
10	0.00110011
100	
1000	
-101	
11	
110	
-101	
1	
10	
100	
1000	
-101	
11	
110	
-101	
1	

What we see is that we get a cycle: every 4 Euclidean divisions, we get a quotient $2r0011$ and a remainder 1. That means that we need an infinite series of this bit pattern 0011 to represent $1/5$ in base 2. Let us see how Pharo dealt to convert $(1/5)$ to a Float:

```
(1/5) asFloat significandAsInteger printStringBase: 2.
returns '110011001100110011001100110011001100110011001100110011010'
```

```
(1/5) asFloat exponent.
returns -3
```

That's the bit pattern we expected, except the last bits 001 have been rounded to upper 010. This is the default rounding mode of Float, round to nearest even. We now know why 0.2 is represented inexactly in machine. It's the same mantissa for 0.1, and its exponent is -4 .

So, when we entered $0.1 + 0.2$, we didn't get exactly $(1/10) + (1/5)$. Instead of that we got:

```
0.1 asTrueFraction + 0.2 asTrueFraction.
returns (10808639105689191/36028797018963968)
```

But that's not all the story... Let us inspect the above fraction: Stéf ► *explain what is the lowBit and highBit* ◀

```
10808639105689191 printStringBase: 2.
returns '100110011001100110011001100110011001100110011001100110011001100111'
```



```
(2.8 asTrueFraction roundTo: 0.01 asTrueFraction) asFloat
returns 2.8000000000000003
```

Using 0.01s2 rather than 0.01 let this example appear to work:

```
2.80 truncateTo: 0.01s2
returns 2.80s2
```

```
2.80 roundTo: 0.01s2
returns 2.80s2
```

But it's just a case of luck, the fact that 2.8 is inexact is enough to cause other surprises as illustrated below:

```
2.8 truncateTo: 0.001s3.
returns 2.799s3
```

```
2.8 < 2.800s3.
returns true
```

Truncating in the Float world is absolutely unsafe. Though, using a ScaledDecimal for rounding is unlikely to cause such discrepancy, except when playing with last digits.

1.5 Fun with Inexact representations

To add a nail to the coffin, let's play a bit more with inexact representations. Let us try to see the difference between different numbers:

```
{
((2.8 asTrueFraction roundTo: 0.01 asTrueFraction) - (2.8 predecessor)) abs -> -1.
((2.8 asTrueFraction roundTo: 0.01 asTrueFraction) - (2.8)) abs -> 0.
((2.8 asTrueFraction roundTo: 0.01 asTrueFraction) - (2.8 successor)) abs -> 1.
} detectMin: [:e | e key ]

returns the pair
0.0->1
```

The following expression returns 0.0->1, which means that: (2.8 asTrueFraction roundTo: 0.01 asTrueFraction) asFloat = (2.8 successor)

But remember that

```
(2.8 asTrueFraction roundTo: 0.01 asTrueFraction) ~= (2.8 successor)
```

It must be interpreted as the nearest Float to (2.8 asTrueFraction roundTo: 0.01 asTrueFraction) is (2.8 successor).

If you want to know how far it is, then get an idea with:

```
((2.8 asTrueFraction roundTo: 0.01 asTrueFraction) - (2.8 successor asTrueFraction))  
  asFloat  
returns -2.0816681711721685e-16
```

1.6 Conclusion

Floats are inexact numbers. Pay attention when you manipulate them. There are much more things to know about floats, and if you are advanced enough, it would be a good idea to check this link from the wikipedia page "What Every Computer Scientist Should Know About Floating-Point Arithmetic".