

Chapter 1

Fun with Floats

with the participation of:

Nicolas Cellier (nicolas.cellier.aka.nice@gmail.com)

Floats are inexact by nature and this can confuse programmers. This chapter introduces this problem and presents some practical solutions to it. The basic message is that Floats are what they are: inexact but fast numbers.

Note that most of the situations described in this chapters are consequences on how Floats are structured by the hardware and are not tied to Pharo. The very same problems in others programming languages.

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  
→ 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 influencing our understanding. Some approaches print a simpler representation of reality than others. In early versions of Pharo printing $0.1 + 0.2$ were printing 0.3 , now it prints 0.30000000000000004 . This change was guided by the idea that it is better not to lie to the user. Showing the inexactness of a float is better than hiding it because one day or another we can be deeply bitten by them.

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

0.3 printString
  → '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
  → '3FD3333333333334'

0.3 hex
  → '3FD3333333333333'
```

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

```
(0.1 + 0.2) storeString
  → '0.30000000000000004'

0.3 storeString
  → '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
  → true

0.3 closeTo: (0.1 + 0.2)
  → true
```

The method `closeTo:` verify that the two compared numbers have less than 0.0001 of difference. Here is its source code.

```
closeTo: num
  "are these two numbers close?"
  num isNumber ifFalse: [^self = num] ifError: [false].
  self = 0.0 ifTrue: [^num abs < 0.0001].
  num = 0 ifTrue: [^self abs < 0.0001].
  ^self = num asFloat
  or: [(self - num) abs / (self abs max: num abs) < 0.0001]
```

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


```
Float precision.
→ 53
```

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

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

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

Until there we've 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. For example, 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).
→ false
```

```
(1/5) = 0.2
→ false
```

Let us see in detail how we could get the fractional bits of $1/5$ *i.e.*, $2r1/2r101$. For that, we must lay out the division:

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

The denominator is a power of 2 as we expect, but we need 54 bits of precision to store the numerator... Float only provides 53. There will be another rounding error to fit into Float representation:

```
(0.1 asTrueFraction + 0.2 asTrueFraction) asFloat = (0.1 asTrueFraction + 0.2
  asTrueFraction).
→ false

(0.1 asTrueFraction + 0.2 asTrueFraction) asFloat significandAsInteger.
→ '10011001100110011001100110011001100110011001100110011001100110100'
```

To summarize what happened, including conversions of decimal representation to Float representation:

```
(1/10) asFloat      0.1  inexact (rounded to upper)
(1/5) asFloat      0.2  inexact (rounded to upper)
(0.1 + 0.2) asFloat ... inexact (rounded to upper)
```

3 inexact operations occurred, and, bad luck, the 3 rounding operations were all to upper, thus they did cumulate rather than annihilate. On the other side, interpreting 0.3 is causing a single rounding error $(3/10)$ asFloat. We now understand why we cannot expect $0.1 + 0.2 = 0.3$.

As an exercise, you could show why $1.3 * 1.3 \neq 1.69$.

1.3 With floats, printing is inexact

One of the biggest trap we learned with above example is that despite the fact that 0.1 is printed '0.1' as if it were exact, it's not. The name `absPrintExactlyOn:base:` used internally by `printString` is a bit confusing, it does not print exactly, but it prints the shortest decimal representation that will be rounded to the same Float when read back (Pharo always converts the decimal representation to the nearest Float).

Another message exists to print exactly, you need to use `printShowingDecimalPlaces:` instead. As every finite Float is represented internally as a Fraction with a denominator being a power of 2, every finite Float has a decimal representation with a finite number of decimals digits (just multiply numerator and denominator with adequate power of 5, and you'll get the digits). Here you go:

```
0.1 asTrueFraction denominator highBit.
→ 56
```

This means that the fraction denominator is 2^{55} and that you need 55 decimal digits after the decimal point to really print internal representation of 0.1 exactly.

```
0.1 printShowingDecimalPlaces: 55.
→ '0.1000000000000000055511151231257827021181583404541015625'
```

And you can retrieve the digits with:

```
0.1 asTrueFraction numerator * (5 raisedTo: 55).
→ 1000000000000000055511151231257827021181583404541015625
```

You can just check our result with:

```
1000000000000000055511151231257827021181583404541015625/(10 raisedTo: 55) =
0.1 asTrueFraction
→ true
```

You see that printing the exact representation of what is implemented in machine would be possible but would be cumbersome. Try printing $1.0e-100$ exactly if not convinced.

1.4 Float rounding is also inexact

While float equality is known to be evil, you have to pay attention to other aspects of floats. Let us illustrate that point with the following example.

```
2.8 truncateTo: 0.01
→ 2.8000000000000003

2.8 roundTo: 0.01
→ 2.8000000000000003
```

It is surprising but not false that `2.8 truncateTo: 0.01` does not return `2.8` but `2.8000000000000003`. This is because `truncateTo:` and `roundTo:` perform several operations on floats: inexact operations on inexact numbers can lead to cumulative rounding errors as you saw above, and that's just what happens again.

Even if you perform the operations exactly and then round to nearest Float, the result is inexact because of the initial inexact representation of `2.8` and `0.01`.

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

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

```
2.80 truncateTo: 0.01s2
→ 2.80s2
```

```
2.80 roundTo: 0.01s2
  → 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.
  → 2.799s3
```

```
2.8 < 2.800s3.
  → 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 ]

  → 0.0->1
```

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

But remember that

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

It must be interpreted as the nearest Float to $(2.8 \text{ asTrueFraction roundTo: } 0.01 \text{ asTrueFraction})$ 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
  → -2.08166817111721685e-16
```

1.6 Chapter summary

Floats are approximation of a real number by being able to support a wide range of decimal values. This chapters has reviewed the following points

- Never use `=` to compare floats (*e.g.*, `(0.1 + 0.2) = 0.3` returns false)
- Use `closeTo`: instead (*e.g.*, `(0.1 + 0.2) closeTo: 0.3` returns true)
- A float number is represented in base as *sign* \times *mantissa* $\times 2^{\text{exponent}}$ (*e.g.*, $1.2345 = 12345 \times 10^{-4}$)
- `truncateTo`: and `roundTo`: do not always work when truncating or rounding up float (*e.g.*, `2.8 roundTo: 0.01` returns `2.800...003`)

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" (<http://www.validlab.com/goldberg/paper.pdf>).