

Review. Possibilities for binary representations of real numbers:

- fixed-point number: $\pm x.y$ with x and y of a certain number of bits
- floating-point number: $\pm 1.x \cdot 2^y$ with x and y of a certain number of bits
 - IEEE 754, single precision: 32 bit (1 bit for sign, 23 bit for significand, 8 bit for exponent)
 - IEEE 754, double precision: 64 bit (1 bit for sign, 52 bit for significand, 11 bit for exponent)

Example 11. (reasons for floats) Almost universally, major programming languages use floating-point numbers for representing real numbers. Why not fixed-point numbers?

Solution. Fixed-point numbers have some serious issues for scientific computation. Most notably:

- Scaling a number typically results in a loss of precision.
For instance, dividing a number by 2^r and then multiplying it with 2^r loses r digits of precision (in particular, this means that it is computationally dangerous to change units). Make sure that you see that this does not happen for floating-point numbers.
- The range of numbers is limited.
For instance, the largest number is on the order of 2^N where N is the number of bits used for the integer part. On the other hand, a floating-point number can be of the order of 2^{2^M} where M is the number of bits used for the exponent. (Make sure you see how enormous of a difference this is!)

Moreover, as noted in the box below, fixed-point numbers do not really offer anything that isn't already provided by integers. This is the reason why most programming languages don't even offer built-in fixed-point numbers.

Fixed-point numbers are essentially like integers.
For instance, instead of 21.013 (say, seconds) we just work with 21013 (which now is in milliseconds).

Example 12. Give an example where one should not use floats.

Solution. Most notably, one should not use floats when dealing with money. That is because, as we saw earlier, an amount such as 0.10 dollars cannot be represented exactly using a float (when using base 2, as is the default in most programming languages such as Python) and thus will get rounded. This is very problematic when working with money.

Comment. For most purposes, the easiest way to avoid these issues is to store dollar amounts as cents. For the latter we can then simply use integers and work with exact numbers (no rounding).

Recall that a floating-point number (with base 2) is of the form $\pm 1.x \cdot 2^y$ where $1.x$ is the significand and y the exponent. IEEE 754 offers the following choices:

- **single precision:** 32 bit (1 bit for sign, 23 bit for significand, 8 bit for exponent)
- **double precision:** 64 bit (1 bit for sign, 52 bit for significand, 11 bit for exponent)

In IEEE 754, a constant, called a **bias**, is added to the exponent so that all exponents are positive (this avoids using a sign bit for the exponent). Namely, one stores $x + \text{bias}$ where $\text{bias} = 2^7 - 1 = 127$ for single precision ($\text{bias} = 2^{10} - 1$ for double precision).

Example 13. Represent 4.5 as a single precision floating-point number according to IEEE 754.

Solution. $4.5 = 1.125 \cdot 2^2 = \boxed{+} \underbrace{1.001}_{\text{binary}} \cdot 2^2$

The exponent 2 gets stored as $2 + 127 = \boxed{1000,0001}$.

Overall, 4.5 is stored as $\boxed{0} \boxed{1000,0001} \boxed{0010,0000,0000,0000,0000,000}$.

Example 14. Represent -0.1 as a single precision floating-point number according to IEEE 754.

Solution. In a previous example, we computed that $0.1 = (0.0001100110011\dots)_2$.

Hence: $-0.1 = \boxed{-} \underbrace{1.1001,1001\dots}_{\text{binary}} \cdot 2^{-4}$

The exponent -4 gets stored as $-4 + 127 = \boxed{0111,1011}$.

Overall, -0.1 is stored as $\boxed{1} \boxed{0111,1011} \boxed{1001,1001,1001,\dots}$.

Caution. Note that we are not able to store -0.1 exactly. Therefore we need to be careful about how to choose the final bit to best approximate -0.1 . According to IEEE 754, the final bit should be 1 (rather than 0 which we would get if we simply truncated) so that -0.1 gets stored as $\boxed{1} \boxed{0111,1011} \boxed{1001,1001,1001,1001,101}$.

Note that it certainly makes sense to round up the final 0 to 1 because it is followed by 1... (this is similar to us rounding up a final 0 in decimal to 1 if it is followed by 5...).

Example 15. `Python` Explain the following floating-point rounding issue:

```
>>> 0.1 + 0.1 + 0.1 == 0.3
```

```
False
```

```
>>> 0.1 + 0.1 + 0.1
```

```
0.30000000000000004
```

Solution. As we saw in the previous example, 0.1 cannot be stored exactly as a floating-point number (when using base 2). Instead, it gets rounded up slightly. After adding three copies of this number, the error has increased to the point that it becomes visible as in the above output.

IMPORTANT. In the Python code above, we used the operator `==` (two equal signs) to compare two quantities. Note that we cannot use `=` (single equal sign) because that operator is used for assignment (`x = y` assigns the value of `y` to `x`, whereas `x == y` checks whether `x` and `y` are equal).

Comment. As the above issue shows, we should never test two floats x and y for equality. Instead, one typically tests whether the difference $|x - y|$ is less than a certain appropriate threshold. An alternative practical way is to round the floats before comparison (below, we round to 8 decimal digits):

```
>>> round(0.1 + 0.1 + 0.1, 8) == round(0.3, 8)
```

```
True
```

Example 16. `Python` Recall that, to express, say, 0.1 in binary, we compute:

- $2 \cdot 0.1 = \boxed{0}.2$
- $2 \cdot 0.2 = \boxed{0}.4$
- $2 \cdot 0.4 = \boxed{0}.8$
- $2 \cdot 0.8 = \boxed{1}.6$
- $2 \cdot 0.6 = \boxed{1}.2$
- and so on...

The above 5 multiplications with 2 reveal 5 digits after the “decimal” point: $0.1 = (0.00011\dots)_2$. (We can further see that the last four digits repeat; but we will ignore that fact here.)

Let us use Python to do this computation for us. We will start with very basic and naive code, and then upgrade it next time.

```
>>> x = 0.1 # or any value < 1
```

Comment. Everything after the # symbol is considered a comment. This is useful for reminding ourselves of things related to the surrounding code. Comments are usually on a separate line but can be used as above (here, we remind ourselves that the code that follows is not going to handle a number like 2.1 correctly).

To have Python do the above computation for us, we plan to multiply x by 2 (call the result x again), collect a digit (we get that digit as the integer part of x), then subtract that digit from x and repeat. Python has a function called `trunc` which “truncates” a float to its integer part but we need to import it from a package called `math` to make it available.

```
>>> from math import trunc
```

Advanced comment. We can also use `*` in place of `trunc` to import all the functions from the `math` package. However, it is good practice to be explicit about what we need from a package. Note that the function `trunc` is very close to the function `floor` (which computes $\lfloor x \rfloor$, the floor of x , which is the closest integer when rounding down) which also seems appropriate here; however, `floor` returns a float rather than an integer, and we prefer the latter. Also note that we could use `int` (this is a general function that converts an input to an integer) instead of `trunc`. We chose `trunc` because it is more explicitly what we want, and because it gives us a chance to see how to import functions from a package.

We are now ready to compute the first digit:

```
>>> x = 2*x
      digit = trunc(x)
      print(digit)
      x = x-digit
```

0

By copying-and-pasting these four lines four more times, we can produce the next four digits:

```
>>> x = 2*x
      digit = trunc(x)
      print(digit)
      x = x-digit
      x = 2*x
      digit = trunc(x)
      print(digit)
      x = x-digit
      x = 2*x
      digit = trunc(x)
      print(digit)
      x = x-digit
      x = 2*x
      digit = trunc(x)
      print(digit)
      x = x-digit
```

0

0

1

1

Clearly, we should not have to copy-and-paste repeated code like this. We will fix this issue next time, as well as discuss several other important improvements.