

Intro to Research Computing with Python: Numerics

Erik Spence

SciNet HPC Consortium

6 November 2014

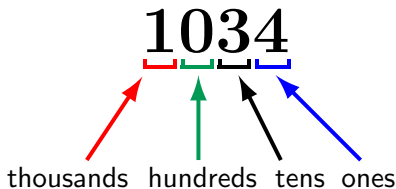
Today's class

Today we will discuss the following topics:

- Numbers. How are they represented and why.
- How computers store different types of numbers.
- The kinds of errors can creep into your calculations, if you're not careful.

How do we represent quantities?

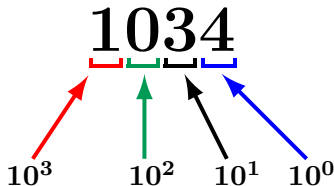
- We use numbers, of course.
- In grade school we are taught that numbers are organized in columns of digits. We learn the names of these columns.
- The numbers are understood as multiplying the digit in the column by the number that names the column.



$$1034 = (1 \times 1000) + (0 \times 100) + (3 \times 10) + (4 \times 1)$$

Other ways to represent a quantity

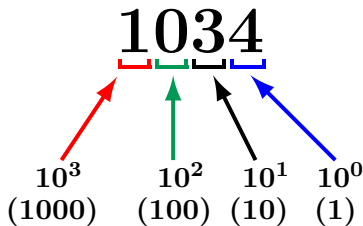
- Instead of using 'tens' and 'hundreds', let's represent the columns by powers of what we will call the 'base'.
- Our normal way of representing numbers is 'base 10', also called decimal.
- Each column represents a power of ten, and the coefficient can be one of 10 numerals (0-9).



$$1034 = (1 \times 10^3) + (0 \times 10^2) + (3 \times 10^1) + (4 \times 10^0)$$

You can choose any base you want

How do we represent the quantity 1034 if we change bases? What about base 7 (septenary)?

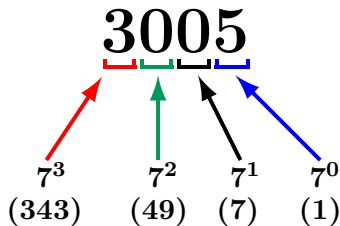
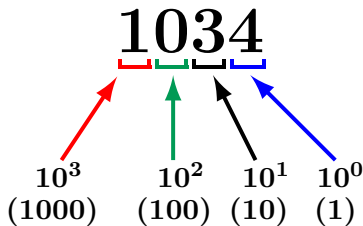


$$1034 = (1 \times 10^3) + (0 \times 10^2) + (3 \times 10^1) + (4 \times 10^0)$$

In base 7 the numerals have the range 0-6.

You can choose any base you want

How do we represent the quantity 1034 if we change bases? What about base 7 (septenary)?



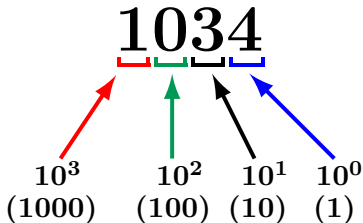
$$1034 = (1 \times 10^3) + (0 \times 10^2) + (3 \times 10^1) + (4 \times 10^0)$$

$$1034 = (3 \times 7^3) + (0 \times 7^2) + (0 \times 7^1) + (5 \times 7^0)$$

In base 7 the numerals have the range 0-6.

Who cares?

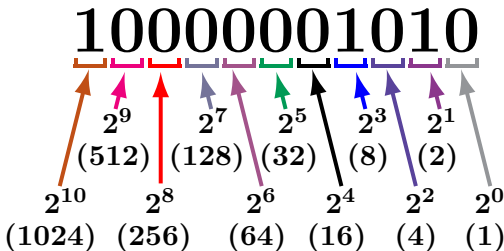
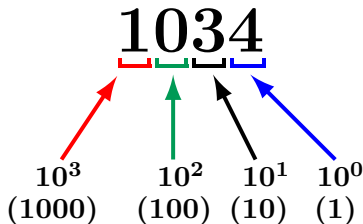
The reason we care is because computers do not use base 10 to store their data. Computers use base 2 (binary). The numerals have the range 0-1.



$$1034 = (1 \times 10^3) + (0 \times 10^2) + (3 \times 10^1) + (4 \times 10^0)$$

Who cares?

The reason we care is because computers do not use base 10 to store their data. Computers use base 2 (binary). The numerals have the range 0-1.



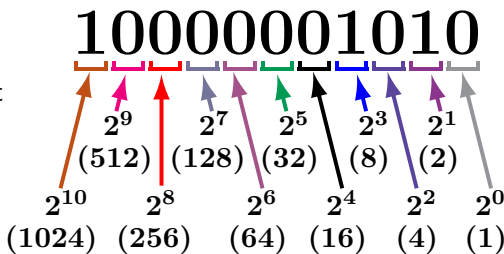
$$1034 = (1 \times 10^3) + (0 \times 10^2) + (3 \times 10^1) + (4 \times 10^0)$$

$$\begin{aligned} 1034 = & (1 \times 2^{10}) + (0 \times 2^9) + (0 \times 2^8) + (0 \times 2^7) \\ & + (0 \times 2^6) + (0 \times 2^5) + (0 \times 2^4) + (1 \times 2^3) \\ & + (0 \times 2^2) + (1 \times 2^1) + (0 \times 2^0) \end{aligned}$$

Why do computers use binary numbers?

Why use binary?

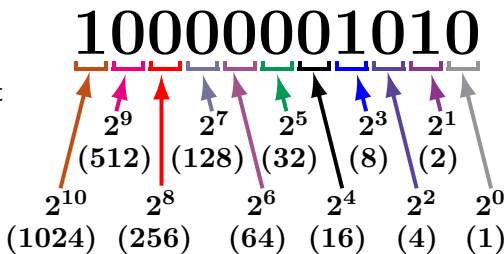
- Modern computers operate using circuits that have one of two states: 'on' or 'off'.
- This choice is related to the complexity and cost of building binary versus ternary circuitry.
- Binary numbers are like series of 'switches': each digit is either 'on' or 'off'.
- Each 'switch' in the number is called a 'bit'.



Why do computers use binary numbers?

Why use binary?

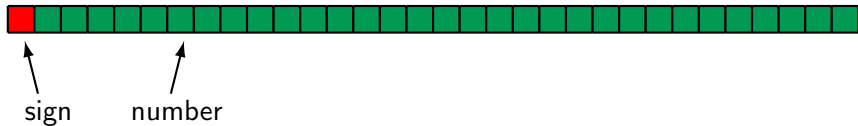
- Modern computers operate using circuits that have one of two states: 'on' or 'off'.
- This choice is related to the complexity and cost of building binary versus ternary circuitry.
- Binary numbers are like series of 'switches': each digit is either 'on' or 'off'.
- Each 'switch' in the number is called a 'bit'.



Pretend that each finger on one of your hands represents one bit. Count to 31 ($2^5 - 1$) on one hand in binary!

Integers

- All integers are exactly representable.
- Different sizes of integer variables are available, depending on your hardware, OS, and programming language.
- For most languages, a typical integer is 32 bits, 1 bit for the sign.
- Finite range: can go from -2^{31} to $2^{31} - 1$ (-2,147,483,648 to 2,147,483,647).
- Unsigned integers: $0 \dots 2^{32} - 1$.
- All operations (+, -, *) between representable integers are represented unless there is overflow.



A typical int = 32 bits = 4 bytes.

Long integers

- Long integers are like regular integers, just with a bigger memory size, usually 64 bits.
- And consequently a bigger range of numbers.
- One bit for sign.
- can go from -2^{63} to $2^{63} - 1$
- -9,223,372,036,854,775,808 to 9,223,372,036,854,775,807.
- Unsigned long integers: $0 \dots 2^{64} - 1$.



A typical long int = 64 bits = 8 bytes.

Integers in Python

Python offers two default types of integers:

- "plain integers":
 - ▶ All integers are plain by default unless they are too big.
 - ▶ These are implemented using long integers in C. This gives them, depending on the system, at least 32 bits of precision.
 - ▶ The maximum value can be found by checking the `sys.maxint` value.
- "long integers":
 - ▶ Have infinite precision.
 - ▶ Are invoked using the `long(something)` function, or by placing an "L" after the number.

```
>>> import sys; print sys.maxint
9223372036854775807
>>> a = 10; b = 10L
>>> type(a)
int
>>> type(b)
long
```

Fixed point numbers

How do we deal with decimal places?

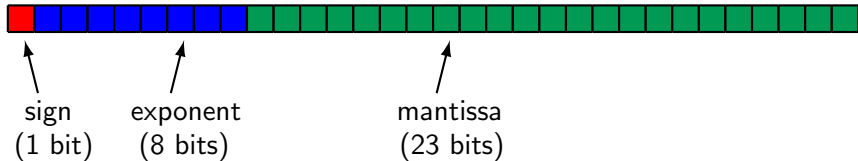
- We could treat real numbers like integers: 0 ... INT_MAX, and only keep, say, the last two digits behind the decimal point.
- This is known as 'fixed point' numbers, since the decimal place is always in the same spot.
- It is often used for financial timeseries data, since they only use a finite number of decimal places.
- But this is terrible for scientific computing. Relative precision varies with magnitude; we need to be able to represent small and large numbers at the same time.
- If you want to deal with fixed point numbers, look into the "decimal" package.

Floating point numbers

- Analog of numbers in scientific notation.
- Inclusion of an exponent means the decimal point is 'floating'.
- Again, one bit is dedicated to sign.

-1.34×10^{-7}

sign mantissa base exponent



A typical single precision real = 32 bits = 4 bytes.

A typical double precision real = 64 bits = 8 bytes.

Floats in Python

Python offers two types of floating point numbers:

- "floating point numbers":
 - ▶ Based on the C double type.
 - ▶ You can specify the exponent by putting "e" in your number.
 - ▶ Information about floats on your system can be found in `sys.float_info`.
- "complex numbers":
 - ▶ Have a real and imaginary part, both of which are floats.
 - ▶ Use `z.real` and `z.imag` to access individual parts.

```
>>> import sys; print sys.float_info
sys.floatinfo(max=1.7976931348623157e+308, max_exp=1024, max_10_exp=308,
min=2.2250738585072014e-308, min_exp=-1021, min_10_exp=-307, dig=15,
mant_dig=53, epsilon=2.2204460492503131e-16, radix=2, rounds=1)
>>> a = complex(1.,3.0); print a
(1+3j)
>>> b = 1.0 + 2.j; print b.imag
2.0
```


Special “numbers”

This format for storing floating point numbers comes from the IEEE 754 standard.

There's room in the format for the storing of a few special numbers.

- Signed infinities ($+\text{Inf}$, $-\text{Inf}$): result of overflow, or divide by zero.
- Signed zeros: signed underflow, or divide by ± 0 .
- Not a Number (NaN): square root of a negative number, $0/0$, Inf/Inf , *etc.*
- The events which lead to these are usually errors, and can be made to cause exceptions.

Errors in floating point mathematics

There are errors inherent in using finite-length floating point variables.

- Except for numbers that fit exactly into a base two representation, assigning a real number to a floating point variable involves truncation.
- Think about how you represent $1/3$. Is it 0.3? 0.33? 0.333?
- You end up with an error of $1/2$ ULP (Unit in Last Place).

```
In [1]: a = 0.1
```

```
In [2]: print a  
Out[2]: 0.1
```

```
In [3]: a  
Out[3]: 0.10000000000000001
```

In base two, 0.1 is an infinitely repeating fraction:

0.0001100110011001100110011...

Single precision: 1 part in $2^{-24} \sim 6\text{e-}8$.

Double precision: 1 part in $2^{-53} \sim 1\text{e-}16$.

Testing for equality

Never ever ever ever test for equality with floating point numbers!

- Because of rounding errors in floating point numbers, you don't know exactly what you're going to get.
- Instead, test to see if the difference is below some tolerance that is near epsilon.
- Testing for equality with integers is ok, however, because integers are exact.

```
In [4]: a = 0.1 * 0.1
```

```
In [5]: b = 0.01
```

```
In [6]: (a == b)
```

```
Out[6]: False
```

```
In [7]: a
```

```
Out[7]: 0.010000000000000002
```

```
In [8]: b
```

```
Out[8]: 0.01
```

```
In [9]: (abs(a - b) < 1e-15)
```

```
Out[9]: True
```

Floating point mathematics

One must be very careful when doing floating point mathematics.

Fire up Python and try the examples on the right.

What went wrong?

```
In [10]: print 1.
```

```
Out[10]: 1.0
```

```
In [11]: print 1.e-18
```

```
Out[11]: 1e-18
```

```
In [12]: print (1. - 1.) + 1.e-18
```

```
In [13]: print (1. + 1.e-18) - 1.
```

```
In [14]: print 1. + 1.e-18
```

Floating point mathematics

One must be very careful when doing floating point mathematics.

Fire up Python and try the examples on the right.

What went wrong?

```
In [10]: print 1.
```

```
Out[10]: 1.0
```

```
In [11]: print 1.e-18
```

```
Out[11]: 1e-18
```

```
In [12]: print (1. - 1.) + 1.e-18
```

```
Out[12]: 1e-18
```

```
In [13]: print (1. + 1.e-18) - 1.
```

```
In [14]: print 1. + 1.e-18
```

Floating point mathematics

One must be very careful when doing floating point mathematics.

Fire up Python and try the examples on the right.

What went wrong?

```
In [10]: print 1.
```

```
Out[10]: 1.0
```

```
In [11]: print 1.e-18
```

```
Out[11]: 1e-18
```

```
In [12]: print (1. - 1.) + 1.e-18
```

```
Out[12]: 1e-18
```

```
In [13]: print (1. + 1.e-18) - 1.
```

```
Out[13]: 0.0
```

```
In [14]: print 1. + 1.e-18
```

Floating point mathematics

One must be very careful when doing floating point mathematics.

Fire up Python and try the examples on the right.

What went wrong?

```
In [10]: print 1.
```

```
Out[10]: 1.0
```

```
In [11]: print 1.e-18
```

```
Out[11]: 1e-18
```

```
In [12]: print (1. - 1.) + 1.e-18
```

```
Out[12]: 1e-18
```

```
In [13]: print (1. + 1.e-18) - 1.
```

```
Out[13]: 0.0
```

```
In [14]: print 1. + 1.e-18
```

```
Out[14]: 1.0
```

Machine epsilon

Let's do some addition, to demonstrate what went wrong.

- Problem: $1.0 + 0.001$
- Let's work in base 10.
- Let's assume that we only have a mantissa precision of 3, and exponent precision of 2.

$$\begin{array}{r} 1.00 \times 10^0 \\ + 1.00 \times 10^{-3} \\ \hline \end{array}$$

$$\begin{array}{r} 1.00 \times 10^0 \\ + 0.001 \times 10^0 \\ \hline 1.00 \times 10^0 \end{array}$$

Machine epsilon

Let's do some addition, to demonstrate what went wrong.

- Problem: $1.0 + 0.001$
- Let's work in base 10.
- Let's assume that we only have a mantissa precision of 3, and exponent precision of 2.
- So what happened?
- Mantissa only has a precision of 3! The final answer is beyond the range of the mantissa!

$$\begin{array}{r} 1.00 \times 10^0 \\ + 1.00 \times 10^{-3} \\ \hline \end{array}$$

$$\begin{array}{r} 1.00 \times 10^0 \\ + 0.001 \times 10^0 \\ \hline 1.00 \times 10^0 \end{array}$$

Machine epsilon

Machine epsilon gives you the limits of the precision of the machine.

- Machine epsilon is defined to be the smallest x such that $1 + x \neq 1$.
- (or sometimes, the largest x such that $1 + x = 1$.)
- Machine epsilon is named after the mathematical term for a small positive infinitesimal.

```
In [15]: print 1.
```

```
Out[15]: 1.0
```

```
In [16]: print 1.e-18
```

```
Out[16]: 1e-18
```

```
In [17]: print (1. - 1.) + 1.e-18
```

```
Out[17]: 1e-18
```

```
In [18]: print (1. + 1.e-18) - 1.
```

```
Out[18]: 0.0
```

```
In [19]: print 1. + 1.e-18
```

```
Out[19]: 1.0
```

What's your epsilon?

You can find your approximate machine epsilon by repeatedly halving a number and testing it.

```
# myepsilon.py
def myepsilon():

    # Initialize our epsilon.
    eps = 1.0

    # Is (1 + eps) > 1?
    while ((1. + eps) > 1.):
        # If it is, divide and print it.
        eps = eps / 2.
        # Change the number of digits
        # printed so we can see them
        # all.
        print '%1.8e %1.18f' % \
            (eps, (1. + eps))
```

```
In [20]: import myepsilon
In [21]: myepsilon.myeepsilon()
.
.
.
1.77635684e-15 1.0000000000000001776
8.88178420e-16 1.0000000000000000888
4.44089210e-16 1.0000000000000000444
2.22044605e-16 1.0000000000000000222
1.11022302e-16 1.0000000000000000000
In [22]:
In [22]: import sys
In [23]: sys.float_info.epsilon
2.2204460492503131e-16
```

The epsilon is about $1e-16$ for my desktop, as expected for double precision.

The limits of precision: look out below!

Problems will occur when the result of a calculation spans more orders of magnitude than the inverse of machine epsilon.

Try the following:

- For the range of numbers between 0 and 2, repeatedly take square roots of the numbers, then repeatedly square the numbers.
- Plot the result, from 0..2.
- What should you get? What do you get?
- Loss of precision in early stages of a calculation causes problems.

```
# precision.py
from numpy import sqrt
def sqrts(x):

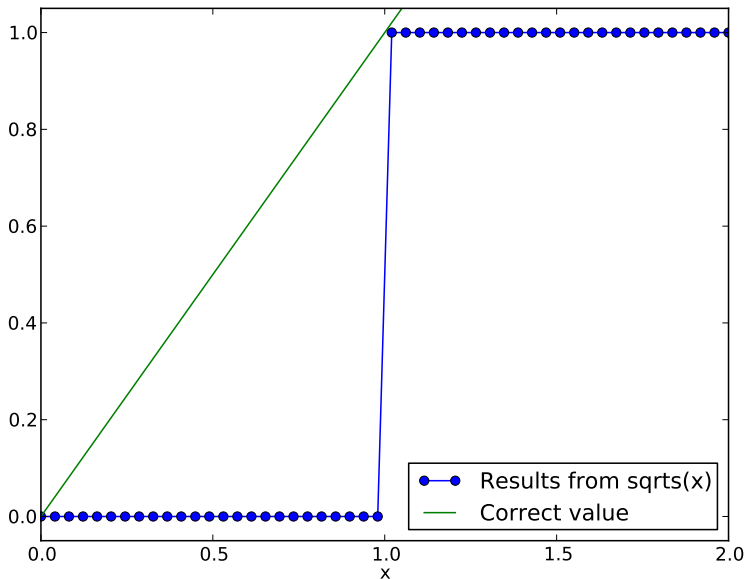
    # Make a copy of the argument.
    y = x

    # Repeatedly sqrt, then square.
    for i in xrange(128):
        y = sqrt(y)
    for i in xrange(128):
        y = y * y

    return y
```

```
In [22]: import precision
In [23]: x = linspace(0., 2., 50)
In [24]: y = precision.sqrts(x)
In [25]: plot(x, y, 'o-')
```

Precision problem: uh oh



Precision problem: what happened?

```
# precision.py
from numpy import sqrt
def sqrts(x):
    y = x
    for i in xrange(128):
        y = sqrt(y)
        print '%1i %1.16f' % (i,y)
    for i in xrange(128):
        y = y * y
        print '%1i %1.16f' % (i,y)
    return y
```

If the argument is below 1.0, sqrt pushes it up to epsilon below 1.0.

If the argument is above 1.0, sqrt pulls it down to exactly 1.0.

```
In [26]: sqrts(0.1)
0 0.3162277660168379
1 0.5623413251903491
.
.
.
126 0.9999999999999999
127 0.9999999999999999
0 0.9999999999999998
1 0.9999999999999996
2 0.9999999999999991
3 0.9999999999999982
.
.
.
126 0.0000000000000000
127 0.0000000000000000
Out[26]: 0.0
```

```
In [27]: sqrts(1.9)
0 1.3784048752090221
1 1.1740548859440185
.
.
.
126 1.0000000000000000
127 1.0000000000000000
0 1.0000000000000000
1 1.0000000000000000
2 1.0000000000000000
3 1.0000000000000000
.
.
.
126 1.0000000000000000
127 1.0000000000000000
Out[27]: 1.0
```

Beware: subtraction

Be very wary of subtracting very similar numbers.

- Problem: subtract 1.22 from 1.23.
- Assume that we only have a mantissa precision of 3, and exponent precision of 2.
- By performing this subtraction, we eliminate most of the information, and end up with 'catastrophic cancellation'.
- We go from 3 significant digits to 1.
- Dangerous in intermediate results.

3 sig. digits

$$\begin{array}{r} 1.23 \times 10^0 \\ - 1.22 \times 10^0 \\ \hline 1.00 \times 10^{-2} \end{array}$$

1 sig. digit

The same problem can occur when dividing large numbers.

Overflow

Overflow occurs when the result of a calculation exceeds the memory size of the variable.

- 8-bit integers have a range of -128 to 127.

```
In [27]: from numpy import int8
```

```
In [28]: a = int8(10)
```

```
In [29]: a
```

```
Out[29]: 10
```

```
In [30]: a.dtype
```

```
Out[30]: dtype('int8')
```

```
In [31]: type(a)
```

```
Out[31]: numpy.int8
```

```
In [32]: a * a
```

```
Out[32]: 100
```

```
In [33]: a * a * a
```

```
Out[33]: -24
```

```
In [34]:
```


Overflow

Overflow occurs when the result of a calculation exceeds the memory size of the variable.

- 8-bit integers have a range of -128 to 127.

```
In [27]: from numpy import int8
```

```
In [28]: a = int8(10)
```

```
In [29]: a
```

```
Out[29]: 10
```

```
In [30]: a.dtype
```

```
Out[30]: dtype('int8')
```

```
In [31]: type(a)
```

```
Out[31]: numpy.int8
```

```
In [32]: a * a
```

```
Out[32]: 100
```

```
In [33]: a * a * a
```

```
Out[33]: -24
```

```
In [34]: int8(1000)
```

```
Out[34]: -24
```

```
In [35]:
```

Overflow

Overflow occurs when the result of a calculation exceeds the memory size of the variable.

- 8-bit integers have a range of -128 to 127.
- When Python calculates a quantity, it up-casts all of the variables to the 'largest' variable type in the calculation.
 - ▶ int are converted to long ints
 - ▶ ints are converted to floats
 - ▶ single precision floats are converted to double.
- Always be sure to use variables that are big enough for what you're doing.

```
In [27]: from numpy import int8
```

```
In [28]: a = int8(10)
```

```
In [29]: a
```

```
Out[29]: 10
```

```
In [30]: a.dtype
```

```
Out[30]: dtype('int8')
```

```
In [31]: type(a)
```

```
Out[31]: numpy.int8
```

```
In [32]: a * a
```

```
Out[32]: 100
```

```
In [33]: a * a * a
```

```
Out[33]: -24
```

```
In [34]: int8(1000)
```

```
Out[34]: -24
```

```
In [35]: a * a * int16(a)
```

```
Out[35]: 1000
```

```
In [36]: a * float(a) * int16(a)
```

```
Out[36]: 1000.0
```

Underflow

An underflow error is the opposite of an overflow error: you are attempting to make a number which is smaller than the variable can hold.

- 32-bit floats integers have a range of $-3.4e38$ to $+3.4e38$
- An overflow error will result if you attempt to go beyond this range.
- An underflow error results if you try to go too small.

```
In [37]: from numpy import float32
In [38]:
In [38]: float32(-1e35)
Out[38]: -1e+35
In [39]: float32(-1e44)
Out[39]: -inf
In [40]:
In [40]: float32(1e-40)
Out[40]: 9.9999461e-41
In [41]: float32(1e-44)
Out[41]: 9.8090893e-45
In [42]: float32(1e-46)
Out[42]: 0.0
```

Summary: things to remember

- Integers are stored exactly.
- Floating point numbers are, in general, NOT stored exactly. Rounding error will cause the number to be slightly off.
- DO NOT test floating point numbers for equality. Instead test $(\text{abs}(a - b) < \text{cutoff})$.
- Know the approximate value of epsilon for the machine that you are using.
- Know the limits of your precision: if your calculations span as many orders of magnitude as the inverse of epsilon you're going to lose precision.
- Try not to subtract floating point numbers that are very close to one another. 'Catastrophic cancellation' leads to loss of precision.
- Be aware of overflow and underflow: use variable sizes that are appropriate for your problem.

Homework 1

- 1 Write a program, called `DecimalToTernary`, which takes as its argument a base-10 integer, less than 6561, and returns an array which contains the argument's ternary (base-3) form.

```
In [37]:
```

```
In [37]: DecimalToTernary(149)
```

```
Out[37]: array([0, 0, 0, 1, 2, 1, 1, 2])
```

```
In [38]:
```

Do NOT use Numpy's "`base_repr`" function.

Homework 1, continued

- 2 Write a program, called CalcOverflow, which, given an argument $m > 1.0$, returns the minimum value of integer n that generates an overflow error when calculating m^n .

Note that Python will throw a runtime error when it encounters an overflow; you must catch this exception:

```
In [40]: m = 5.
In [41]: m**500
Traceback (most recent call last):
File "<stdin>", line 1, in <module>
OverflowError: (34, 'Numerical result out of range')
In [42]: try:
...:     m**500
...: except:
...:     print "eeek!"
...:
eeek!
In [43]:
```

Homework 1, continued

- ③ Write a program, called `CalcUnderflow`, which, given an argument $m > 1.0$, returns the minimum value of integer p that generates an underflow error when calculating m^{-p} .

```
In [40]: from mycode import CalcUnderflow
```

```
In [41]: CalcUnderflow(12.3)
```

```
Out[41]: 297
```

```
In [42]:
```

For those that are worried, the questions will get more interesting in the coming weeks.