

Data Science: Numpy

**CPSC 501: Advanced Programming Techniques
Fall 2022**

Jonathan Hudson, Ph.D
Assistant Professor (Teaching)
Department of Computer Science
University of Calgary

Monday, October 3, 2022

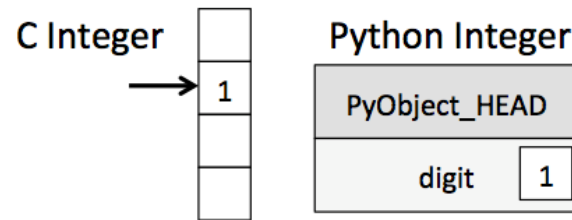


numpy

- **Numerical Python** library (2005 brought together disparate library ideas)
- More efficient data and storage operations as arrays grow larger
- **Python** integer is more than an integer, and a list more than just values

integers?

- **Python** integer is more than an integer
- **Python 3.4** integer

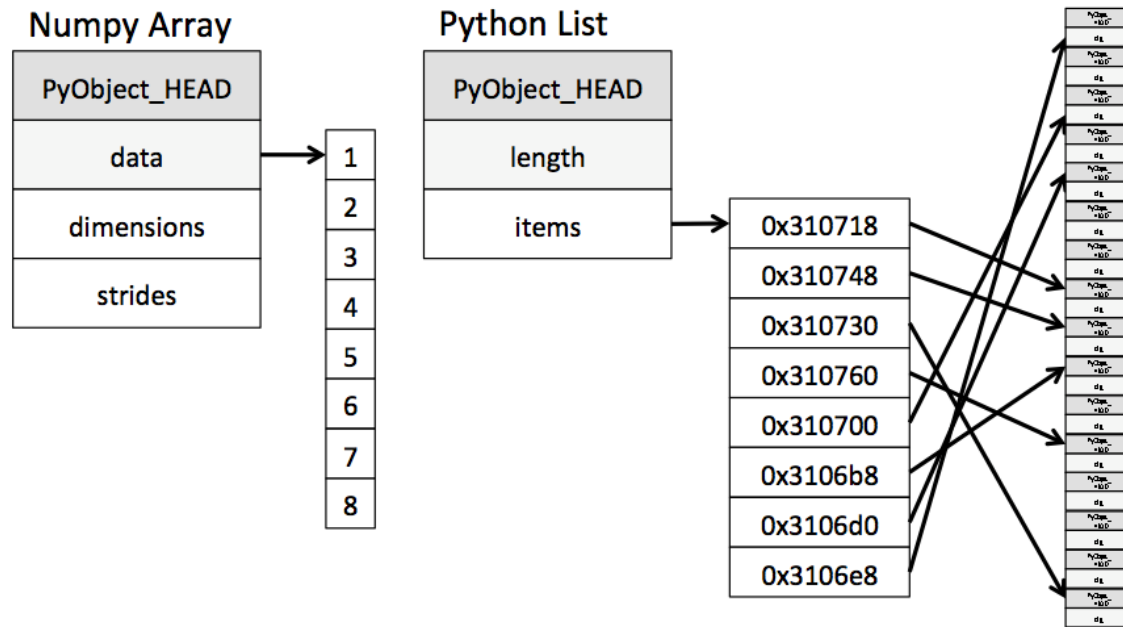


```
1 struct _longobject {  
2     long ob_refcnt;  
3     PyTypeObject *ob_type;  
4     size_t ob_size;  
5     long ob_digit[1];  
6 }
```

- A **C** integer is essentially a label for a position in memory whose bytes encode an integer value.
- A **Python** integer is a pointer to a position in memory containing all the Python object information, including the bytes that contain the integer value.

lists?

- A **Python** list is more than just a list of values



- A **Python** list contains a pointer to a block of pointers, each of which in turn points to a full **Python** object like the **Python** integer we saw earlier.
- **numpy** arrays are a single pointer to a block of contiguous data.

numpy

- **Numerical Python** library
- More efficient data and storage operations as arrays grow larger
- **Python** integer is more than an integer, and a list more than just values
- **numpy** arrays allow us to put data all in one place and drastically improve the our ability to manipulate it quickly.
- In essences what **numpy** does is provide a portal from **Python** through to **C** implementations of storage arrays, allowing us to access the strengths of that language in ways not normally available in **Python**.



```
1 import numpy as np
```

array?

- **Python** does have its own **array** type

```
▶ import array
  L = list(range(10))
  A = array.array("i",L)
  print(A)

↳ array('i', [0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
```

- There are similarities to representation that **numpy** also uses but **numpy** also adds a large range of efficient operations that the **array** type does not have
- Few people use **array** unless there is no way to import or access **numpy**

Create

- Numerous options to create **numpy** arrays

```
import numpy as np
```

```
[5] np.array(range(5))
```

```
array([0, 1, 2, 3, 4])
```

```
[12] np.array([1,2,3,4,5], dtype="float32")
```

```
array([1., 2., 3., 4., 5.], dtype=float32)
```

```
[11] np.zeros(10, dtype=int)
```

```
array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0])
```

```
[13] np.full(3, 10)
```

```
array([10, 10, 10])
```

```
[16] np.random.random(3)
```

```
array([0.65720754, 0.97030252, 0.77293412])
```

Datatypes

- To gain the power of C arrays Python now has to care about all the types that a C programmer has to manage

Data type	Description
<code>bool_</code>	Boolean (True or False) stored as a byte
<code>int_</code>	Default integer type (same as C <code>long</code> ; normally either <code>int64</code> or <code>int32</code>)
<code>intc</code>	Identical to C <code>int</code> (normally <code>int32</code> or <code>int64</code>)
<code>intp</code>	Integer used for indexing (same as C <code>ssize_t</code> ; normally either <code>int32</code> or <code>int64</code>)
<code>int8</code>	Byte (-128 to 127)
<code>int16</code>	Integer (-32768 to 32767)
<code>int32</code>	Integer (-2147483648 to 2147483647)
<code>int64</code>	Integer (-9223372036854775808 to 9223372036854775807)
<code>uint8</code>	Unsigned integer (0 to 255)
<code>uint16</code>	Unsigned integer (0 to 65535)
<code>uint32</code>	Unsigned integer (0 to 4294967295)
<code>uint64</code>	Unsigned integer (0 to 18446744073709551615)
<code>float_</code>	Shorthand for <code>float64</code> .
<code>float16</code>	Half precision float: sign bit, 5 bits exponent, 10 bits mantissa
<code>float32</code>	Single precision float: sign bit, 8 bits exponent, 23 bits mantissa
<code>float64</code>	Double precision float: sign bit, 11 bits exponent, 52 bits mantissa
<code>complex_</code>	Shorthand for <code>complex128</code> .
<code>complex64</code>	Complex number, represented by two 32-bit floats
<code>complex128</code>	Complex number, represented by two 64-bit floats

Speed

- numpy operations speed comes from a large variety of Universal Functions (Ufuncs) which are 'vectorized operations' designed to run at higher speed than if we let Python's loops manage things
- Both computer same answer!

```
[48] import numpy as np
      np.random.seed(0)

      def compute_reciprocals(values):
          output = np.empty(len(values))
          for i in range(len(values)):
              output[i] = 1.0 / values[i]
          return output

      values = np.random.randint(1, 10, size=5)

      print(compute_reciprocals(values))
      print(1.0 / values)
```

```
[0.16666667  1.          0.25          0.25          0.125         ]
[0.16666667  1.          0.25          0.25          0.125         ]
```

Speed

- Both computer same answer!
- Milliseconds versus microseconds (a difference in order of magnitude of 1000 here)

```
big_array = np.random.randint(1, 100, size=100000)
print("Python")
%timeit compute_reciprocals(big_array)
print("numpy")
%timeit (1.0 / big_array)
```

Python

1 loop, best of 5: 225 ms per loop

numpy

The slowest run took 4.95 times longer than the fastest

1000 loops, best of 5: 227 μ s per loop

UFuncs

- Most **Ufuncs** are accessible using straight operators, although we can always use the longer form of `np.divide` for example
- Some don't have operators like `np.abs(array)`
`np.sin(array)`, `np.cos(array)`, etc.
`np.log2(array)`, `np.ln(array)`, etc.
`np.sqrt(array)`, `np.floor(array)`, etc.
- Others available via **scipy** under **special**

```
1 from scipy import special
```

Operator Equivalent ufunc Description

+	<code>np.add</code>	Addition (e.g., <code>1 + 1 = 2</code>)
-	<code>np.subtract</code>	Subtraction (e.g., <code>3 - 2 = 1</code>)
-	<code>np.negative</code>	Unary negation (e.g., <code>-2</code>)
*	<code>np.multiply</code>	Multiplication (e.g., <code>2 * 3 = 6</code>)
/	<code>np.divide</code>	Division (e.g., <code>3 / 2 = 1.5</code>)
//	<code>np.floor_divide</code>	Floor division (e.g., <code>3 // 2 = 1</code>)
**	<code>np.power</code>	Exponentiation (e.g., <code>2 ** 3 = 8</code>)
%	<code>np.mod</code>	Modulus/remainder (e.g., <code>9 % 4 = 1</code>)

Boolean arrays

Operator Equivalent ufunc

<code>==</code>	<code>np.equal</code>	<code>!=</code>	<code>np.not_equal</code>
<code><</code>	<code>np.less</code>	<code><=</code>	<code>np.less_equal</code>
<code>></code>	<code>np.greater</code>	<code>>=</code>	<code>np.greater_equal</code>

Operator Equivalent ufunc

<code>&</code>	<code>np.bitwise_and</code>	<code> </code>	<code>np.bitwise_or</code>
<code>^</code>	<code>np.bitwise_xor</code>	<code>~</code>	<code>np.bitwise_not</code>

- Can use booleans to ask questions, combine booleans with bitwise operators
- Will product boolean arrays which can then be used as 'masks' for filtering

```
x = np.array([5, 2, 3, 4, 1])
print(x < 3)
print(np.any(x<4))
print(np.all(x<4))
print(np.sum(x<4))
print(np.sum((x < 4) & (x > 1)))
print(np.sum((x < 2) | (x > 3)))
print(x[(x < 2) | (x > 3)])
```

```
[False True False False True]
True
False
3
2
3
[5 4 1]
```

Save/Load

- **numpy** has both straight forward save/load, but also ability to save multiple in a form that multiple arrays can be reloaded and referenced by stored key

```
array = np.random.randint(0, 100, size = 15)
np.save('filename', array)
temp = np.load("filename.npy")
print(temp)
array2 = np.random.randint(0, 100, size = 15)
np.savez('filename2', a=array, b=array2)
temp1 = np.load("filename2.npz")
print(temp1['a'])
print(temp1['b'])
```

```
[59 83  3  6 33 96 72 54 36 29 81 85 95 51 65]
[59 83  3  6 33 96 72 54 36 29 81 85 95 51 65]
[27  3 62 15  0 47 56 96 79 40 88 14 92 96 89]
```

Onward to ... pandas.

Jonathan Hudson
jwhudson@ucalgary.ca
<https://pages.cpsc.ucalgary.ca/~jwhudson/>



UNIVERSITY OF
CALGARY