

# NumPy

# Numerical Python

Provides efficient storage and operations on dense data buffers, i.e., arrays.

- ▶ `ndarray` is the fundamental object
- ▶ Vectorized operations on arrays
- ▶ Broadcasting
- ▶ File IO and memory-mapped files

```
1 In [1]: import numpy as np
```

# NumPy Array Element Types

Arrays have elements of homogeneous data type

```
1 In [2]: a = np.array([1, 2, 3.14])
2
3 In [3]: type(a)
4 Out[3]: numpy.ndarray
5
6 In [4]: a
7 Out[4]: array([ 1. ,  2. ,  3.14])
8
9 In [5]: type(a[0])
10 Out[5]: numpy.float64
```

- ▶ Notice that the values were converted to floats.

You can specify an explicit element type with the `dtype` keyword argument:

```
1 In [6]: np.array(nums, dtype='int')
2 Out[6]: array([1, 2, 3])
```

# One-dimensional Arrays

Pass list to `np.array()`:

```
1 In [9]: np.array([1,2,3])
2 Out[9]:
3 array([1, 2, 3])
```

Create a one-dimensional array of zeros, `dtype` defaults to `float`:

```
1 In [10]: np.zeros(4)
2 Out[10]: array([ 0.,  0.,  0.,  0.])
```

`np.arange` similar to Python's built-in `range(start, end, stride)`:

```
1 In [13]: np.arange(0, 10, 2)
2 Out[13]: array([0, 2, 4, 6, 8])
```

# Multi-Dimensional Arrays

Passing nested lists to `np.array()` create multi-dimensional arrays:

```
1 In [9]: np.array([[1,2,3],[4,5,6]])
2 Out[9]:
3 array([[1, 2, 3],
4        [4, 5, 6]])
```

Create a multi-dimensional array of 1s with element type `int`. Note that first argument is a tuple of array dimensions.

```
1 In [11]: np.ones((2, 3), dtype=int)
2 Out[11]:
3 array([[1, 1, 1],
4        [1, 1, 1]])
```

Create a 2-d array of the same element values:

```
1 In [12]: np.full((2, 3), 2.72)
2 Out[12]:
3 array([[ 2.72,  2.72,  2.72],
4        [ 2.72,  2.72,  2.72]])
```

# Creating Arrays of Random Numbers

Create a  $2 \times 3$  array of values uniformly distributed between 0 and 1:

```
1 In [28]: np.random.random((2, 3))
2 Out[28]:
3 array([[ 0.93923457,  0.41299137,  0.07451052],
4        [ 0.32800936,  0.44435825,  0.4520937 ]])
```

Normally distributed with  $\mu = 71.36$  and  $\sigma = 14.79$ :

```
1 In [26]: np.random.normal(71.36, 14.79, (2, 3))
2 Out[26]:
3 array([[ 71.24362489,  61.05019638,  72.25408014],
4        [ 63.03759916,  70.64992342,  75.94207076]])
```

Create a  $2 \times 3$  array of `int` values in the interval  $[1, 11]$ :

```
1 In [29]: np.random.randint(1, 11, (2, 3))
2 Out[29]:
3 array([[9, 8, 6],
4        [9, 5, 9]])
```

3-d identity matrix:

```
1 In [31]: np.identity(3)
2 Out[31]:
3 array([[ 1.,  0.,  0.],
4        [ 0.,  1.,  0.],
5        [ 0.,  0.,  1.]])
```

# NumPy Array Attributes

Given:

```
1 In [33]: a = np.array([[1,2,3], [4,5,6]])
2
3 In [34]: a
4 Out[34]:
5 array([[1, 2, 3],
6        [4, 5, 6]])
```

`ndim` is the number of dimensions:

```
1 In [37]: a.ndim
2 Out[37]: 2
```

`shape` is a tuple giving the number of elements in each dimension:

```
1 In [35]: a.shape
2 Out[35]: (2, 3)
```

`dtype` is the type of the elements

```
1 In [36]: a.dtype
2 Out[36]: dtype('int64')
```

# 1-D Array Indexing and Slicing

1-d arrays similar to Python lists:

```
1 In [41]: a1 = np.arange(10)
2
3 In [44]: a1[1]
4 Out[44]: 1
5
6 In [45]: a1[-1]
7 Out[45]: 9
8
9 In [46]: a1[2:5]
10 Out[46]: array([2, 3, 4])
```

Assignment of single value to a (sub)range /broadcasts/ the value to the (sub)range:

```
1 In [47]: a1[2:5] = 11
2
3 In [48]: a1
4 Out[48]: array([ 0,  1, 11, 11, 11,  5,  6,  7,  8,  9])
```

Notice that the original array is modified.

# 2-D Array Indexing and Slicing

Given:

```
1 In [49]: a3 = np.array([[1,2,3],[4,5,6],[7,8,9]])
2
3 In [50]: a3
4 Out[50]:
5 array([[1, 2, 3],
6        [4, 5, 6],
7        [7, 8, 9]])
```

Single scalar value:

```
1 In [51]: a3[1,1]
2 Out[51]: 5
```

Subarray:

```
1 In [52]: a3[1:, 1:]
2 Out[52]:
3 array([[5, 6],
4        [8, 9]])
```

Single column:

```
1 In [53]: a3[:, 2]
2 Out[53]: array([3, 6, 9])
```



Single row:

# Array Reshaping

## 2-d arrays

```
1 In [62]: a3 = np.arange(1, 13)
2
3 In [63]: a3
4 Out[63]: array([ 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12])
5
6 In [64]: a3.reshape(3, 4)
7 Out[64]:
8 array([[ 1, 2, 3, 4],
9        [ 5, 6, 7, 8],
10       [ 9, 10, 11, 12]])
11
12 In [65]: a3.reshape(4, 3)
13 Out[65]:
14 array([[ 1, 2, 3],
15        [ 4, 5, 6],
16        [ 7, 8, 9],
17        [10, 11, 12]])
```

# Python is slow

- ▶ Consider an array representing pixels of a “one megapixel” image:

```
1 In [20]: image = np.random.randint(0, 256, (1000000, 3))
```

- ▶ This is a deep underwater image which looks very green and we want to increase the “blueness” by 10% [fn:1]. So we write a function to multiply pixel elements by a factor:

```
1 In [60]: def mult_elem(image, n, factor):  
2     ...:     for i in range(len(image)):  
3     ...:         image[i][n] = image[i][n] * factor
```

- ▶ This operation is *slow*:

```
1 In [61]: %timeit mult_elem(image, 2, 1.10)  
2 1.85 s +/- 16.8 ms per loop (mean +/- std. dev. of 7 runs, 1 loop each)
```

- ▶ The equivalent vectorized operation is /300 times faster/:

```
1 In [62]: %timeit image[:, 2] = image[:, 2] * 1.10  
2 6.23 ms +/- .0693 ms per loop (mean +/- std. dev. of 7 runs, 100 loops  
each)
```



[fn:1] I'm not a graphics guy, so just indulge me here.

# Vectorized Operations on Arrays

Operations between compatibly-shaped arrays or between arrays and scalars are *vectorized* – the loop applying the operations to elements of the array(s) is in the compiled C-code layer instead of Python.

```
1 In [114]: np.arange(2, 20, 2) / np.arange(1, 10)
2 Out[114]: array([ 2.,  2.,  2.,  2.,  2.,  2.,  2.,  2.])
```

Smaller array is “broadcast” across the larger array. The simplest example is when the smaller array is a scalar value:

```
1 In [108]: a = np.arange(9)
2
3 In [110]: 2 ## a
4 Out[110]: array([ 1,  2,  4,  8, 16, 32, 64, 128, 256])
5
6 In [111]: 2 ## a.reshape((3, 3))
7 Out[111]:
8 array([[ 1,  2,  4],
9        [ 8, 16, 32],
10       [64, 128, 256]])
```

General broadcasting between multi-dimensional arrays is beyond the scope of this course. See [the NumPy docs](#) for details.



# Masking

First, boolean indexing: you can use a like-shaped array of bools to index into an array, which selects items from the array. The arrays of bools is called a /mask/ and using it to select elements is called /masking/.

```
1 In [175]: xs = np.array([0,1,2,3,4,5,6,7,8,9])
2
3 In [177]: xs[[True, False, True, False, True, False, True, False]]
4 Out[177]: array([0, 2, 4, 6, 8])
```

Since you can create arrays of bools easily with comparison ufuncs, you can combine boolean indexing with broadcasting to easily mask an array:

```
1 In [179]: xs[(xs % 2) == 0]
2 Out[179]: array([0, 2, 4, 6, 8])
```

The comparison operation above is a boolean universal function.

# Boolean UFuncs

Broadcast boolean expressions just like arithmetic expressions:

```
1 In [163]: exam1scores = np.loadtxt('exam1grades.txt')
2
3 In [164]: exam1scores
4 Out[164]:
5 array([ 72.,  72.,  50.,  65.,  60.,  73.,  93.,  88.,  97., ...
6          84.,  75.,  88.,  75.,  86.,  49.,  65.,  69.,  87.])
```

How many people “passed”? First, you can apply a comparison operator to an array to get an array of booleans:

```
1 In [165]: exam1scores > 70
2 Out[165]:
3 array([ True, True, False, False, False, True, True, True, True, ...
4          True, True, True, True, True, False, False, False, True],
      dtype=bool)
```

Then you can apply the `np.sum` aggregation function to count the booleans in the resulting array of booleans:

```
1 In [169]: np.sum(exam1scores > 70)
2 Out[169]: 77
```



## Boolean UFuncs

You can also combine comparisons with logical operators. How many Bs?

```
1 In [173]: np.sum((exam1scores >= 80) & (exam1scores < 90))  
2 Out[173]: 27
```

Note the syntax with single & – NumPy uses efficient bitwise logical operators.

# Array Aggregations

```
1 In [117]: np.arange(10).sum()
2 Out[117]: 45
3
4 In [119]: np.array([8,6,7,5,3,0,9]).min()
5 Out[119]: 0
6
7 In [120]: np.array([8,6,7,5,3,0,9]).max()
8 Out[120]: 9
```

## 2-D Aggregations

```
1 In [131]: np.arange(9).reshape(3,3)
2 Out[131]:
3 array([[0, 1, 2],
4        [3, 4, 5],
5        [6, 7, 8]])
```

We can summarize the values of each column,

```
1 In [132]: np.arange(9).reshape(3,3).min(axis=0)
2 Out[132]: array([0, 1, 2])
3
4 In [133]: np.arange(9).reshape(3,3).max(axis=0)
5 Out[133]: array([6, 7, 8])
```

or summarize the values in each row:

```
1 In [134]: np.arange(9).reshape(3,3).min(axis=1)
2 Out[134]: array([0, 3, 6])
3
4 In [135]: np.arange(9).reshape(3,3).max(axis=1)
5 Out[135]: array([2, 5, 8])
```

Note that *axis* here means *dimension to be collapsed*. So axis 0 means we collapse the rows into one array by applying the aggregation function by column.

# Missing Data

Missing array elements represented as `np.nan` values.

```
1 In [86]: xs = np.array([2, 3, 4, np.nan])
2
3 In [87]: np.mean(xs)
4 Out[87]: nan
```

Ways to handle missing values:

- ▶ Manually masking with `np.isnan`

```
1 In [90]: np.mean(xs[[not np.isnan(x) for x in xs]])
2 Out[90]: 3.0
```

- ▶ Masking using the `numpy.ma` module.

```
1 In [92]: np.ma.masked_invalid(xs).mean()
2 Out[92]: 3.0
```

- ▶ Using NaN-ignoring aggregates:

```
1 In [93]: np.nanmean(xs)
2 Out[93]: 3.0
```

Pandas gives you a few more options, but these cover many cases that come up in practice.

```
np.where
```

`np.where(cond, true_result, false_result)` is a vectorized version of Python's ternary if-else expression.

Here, we double all the even numbers:

```
1 In [12]: a = np.array([[1,2,3], [4,5,6], [7,8,9]])
2
3 In [14]: a
4 Out[14]:
5 array([[1, 2, 3],
6        [4, 5, 6],
7        [7, 8, 9]])
8
9 In [15]: np.where((a % 2) == 0, a * 2, a)
10 Out[15]:
11 array([[ 1,  4,  3],
12        [ 8,  5, 12],
13        [ 7, 16,  9]])
```

Exercise: do that operation above using basic Python on a list of lists.

# Closing Thoughts

Key ideas of NumPy:

- ▶ In-memory arrays of elements with the same data type
- ▶ Static typing of arrays together with vectorized operations of universal functions provide dramatic speed up over equivalent Python code
- ▶ Ufuncs combined with boolean masks makes it easy to partition data
- ▶ Aggregate functions make it easy to summarize data

NumPy is the foundation of the SciPy stack. Even when we don't use it directly (which we often will), it's there underneath the hood.