# Introduction to data science

Patrick Shafto

Department of Math and Computer Science

# Plan for today

- Review rubric and grading

- Numpy and Scipy

- Example of homework

- Rubric for Weds

# HW 2

- Write a short tutorial on dicts (yes, again), functions, and classes

- Cover basic functionality with worked examples

- Due Sunday by 11:59 pm


- Evaluations due Tuesday by 11:59pm

  - Use the rubric!

  - Justify your scores on each section of the rubric!

- Criteria 1:  19/24

  - Simple examples. Easy to understand. Need to add more examples for dictionary and class. More application based examples needs to be added other then simple add and subtract operations. Please add references for the work done.

  - Overall -3-Interpretation,3-Representation,3-Analysis,3-Calculation,3-Assumptions,3-Communications

- Criteria 1:  21/24

- <u>Interpretation</u> - 4 (N/A) <u>Representation</u> - 4 (the examples presented clearly explained the purpose of the code) <u>Calculation</u> - 4 (N/A)

- <u>Application</u> - 2 (Classes were not covered in detail. The tutorial provided very little understanding of the topic) <u>Assumptions</u> - 4 (N/A) <u>Communication</u> - 3 (too much explanation with too little information. Explanations could have been short & crisp)

# HW 3

- Write a short tutorial on numpy and scipy

- Cover basic functionality with worked examples

- Due Sunday by 11:59 pm


- Evaluations due Tuesday by 11:59pm

    - Use the rubric!

    - Justify your scores on each section of the rubric!

# Numpy

# The Basics

NumPy's main object is the homogeneous multidimensional array. It is a table of elements (usually numbers), all of the same type, indexed by a tuple of positive integers. In NumPy dimensions are called *axes*. The number of axes is *rank*.

For example, the coordinates of a point in 3D space `[1, 2, 1]` is an array of rank 1, because it has one axis. That axis has a length of 3. In the example pictured below, the array has rank 2 (it is 2-dimensional). The first dimension (axis) has a length of 2, the second dimension has a length of 3.

```
[[ 1., 0., 0.],
 [ 0., 1., 2.]]
```

NumPy's array class is called `ndarray`. It is also known by the alias `array`. Note that `numpy.array` is not the same as the Standard Python Library class `array.array`, which only handles one-dimensional arrays and offers less functionality. The more important attributes of an `ndarray` object are:

**ndarray.ndim**

> the number of axes (dimensions) of the array. In the Python world, the number of dimensions is referred to as *rank*.

**ndarray.shape**

> the dimensions of the array. This is a tuple of integers indicating the size of the array in each dimension. For a matrix with *n* rows and *m* columns, `shape` will be `(n,m)`. The length of the `shape` tuple is therefore the rank, or number of dimensions, `ndim`.

**ndarray.size**

> the total number of elements of the array. This is equal to the product of the elements of `shape`.

**ndarray.dtype**

> an object describing the type of the elements in the array. One can create or specify dtype's using standard Python types. Additionally NumPy provides types of its own. numpy.int32, numpy.int16, and numpy.float64 are some examples.

## An example

```python
>>>
>>> import numpy as np
>>> a = np.arange(15).reshape(3, 5)
>>> a
array([[ 0,  1,  2,  3,  4],
       [ 5,  6,  7,  8,  9],
       [10, 11, 12, 13, 14]])
>>> a.shape
(3, 5)
>>> a.ndim
2
>>> a.dtype.name
'int64'
>>> a.itemsize
8
>>> a.size
15
>>> type(a)
<type 'numpy.ndarray'>
>>> b = np.array([6, 7, 8])
>>> b
array([6, 7, 8])
>>> type(b)
<type 'numpy.ndarray'>
```

# Array Creation

There are several ways to create arrays.

For example, you can create an array from a regular Python list or tuple using the `array` function. The type of the resulting array is deduced from the type of the elements in the sequences.

```
>>> import numpy as np
>>> a = np.array([2,3,4])
>>> a
array([2, 3, 4])
>>> a.dtype
dtype('int64')
>>> b = np.array([1.2, 3.5, 5.1])
>>> b.dtype
dtype('float64')
```

`array` transforms sequences of sequences into two-dimensional arrays, sequences of sequences of sequences into three-dimensional arrays, and so on.

```
>>>
>>> b = np.array([(1.5,2,3), (4,5,6)])
>>> b
array([[ 1.5,  2. ,  3. ],
       [ 4. ,  5. ,  6. ]])
```

The type of the array can also be explicitly specified at creation time:

```
>>>
>>> c = np.array( [ [1,2], [3,4] ], dtype=complex )
>>> c
array([[ 1.+0.j,  2.+0.j],
       [ 3.+0.j,  4.+0.j]])
```

Often, the elements of an array are originally unknown, but its size is known. Hence, NumPy offers several functions to create arrays with initial placeholder content. These minimize the necessity of growing arrays, an expensive operation.

The function `zeros` creates an array full of zeros, the function `ones` creates an array full of ones, and the function `empty` creates an array whose initial content is random and depends on the state of the memory. By default, the dtype of the created array is `float64`.

```
>>>
>>> np.zeros( (3,4) )
array([[ 0.,  0.,  0.,  0.],
       [ 0.,  0.,  0.,  0.],
       [ 0.,  0.,  0.,  0.]])
>>> np.ones( (2,3,4), dtype=np.int16 )          # dtype can also be specified
array([[[ 1, 1, 1, 1],
        [ 1, 1, 1, 1],
        [ 1, 1, 1, 1]],

       [[ 1, 1, 1, 1],
        [ 1, 1, 1, 1],
        [ 1, 1, 1, 1]]], dtype=int16)
>>> np.empty( (2,3) )                            # uninitialized, output may vary
array([[  3.73603959e-262,   6.02658058e-154,   6.55490914e-260],
       [  5.30498948e-313,   3.14673309e-307,   1.00000000e+000]])
```

To create sequences of numbers, NumPy provides a function analogous to `range` that returns arrays instead of lists.

```
>>>
>>> np.arange( 10, 30, 5 )
array([10, 15, 20, 25])
>>> np.arange( 0, 2, 0.3 )                    # it accepts float arguments
array([ 0. ,  0.3,  0.6,  0.9,  1.2,  1.5,  1.8])
```

When `arange` is used with floating point arguments, it is generally not possible to predict the number of elements obtained, due to the finite floating point precision. For this reason, it is usually better to use the function `linspace` that receives as an argument the number of elements that we want, instead of the step:

```
>>>
>>> from numpy import pi
>>> np.linspace( 0, 2, 9 )                    # 9 numbers from 0 to 2
array([ 0. ,  0.25,  0.5 ,  0.75,  1. ,  1.25,  1.5 ,  1.75,  2. ])
>>> x = np.linspace( 0, 2*pi, 100 )           # useful to evaluate function at lots of points
>>> f = np.sin(x)
```

# Printing Arrays

When you print an array, NumPy displays it in a similar way to nested lists, but with the following layout:

- the last axis is printed from left to right,
- the second-to-last is printed from top to bottom,
- the rest are also printed from top to bottom, with each slice separated from the next by an empty line.

One-dimensional arrays are then printed as rows, bidimensionals as matrices and tridimensionals as lists of matrices.

```
>>>
>>> a = np.arange(6)                        # 1d array
>>> print(a)
[0 1 2 3 4 5]
>>>
>>> b = np.arange(12).reshape(4,3)          # 2d array
>>> print(b)
[[ 0  1  2]
 [ 3  4  5]
 [ 6  7  8]
 [ 9 10 11]]
>>>
>>> c = np.arange(24).reshape(2,3,4)        # 3d array
>>> print(c)
[[[ 0  1  2  3]
  [ 4  5  6  7]
  [ 8  9 10 11]]
 [[12 13 14 15]
  [16 17 18 19]
  [20 21 22 23]]]
```

# Basic Operations

Arithmetic operators on arrays apply *elementwise*. A new array is created and filled with the result.

```
>>> a = np.array( [20,30,40,50] )
>>> b = np.arange( 4 )
>>> b
array([0, 1, 2, 3])
>>> c = a-b
>>> c
array([20, 29, 38, 47])
>>> b**2
array([0, 1, 4, 9])
>>> 10*np.sin(a)
array([ 9.12945251, -9.88031624,  7.4511316 , -2.62374854])
>>> a<35
array([ True, True, False, False], dtype=bool)
```

Unlike in many matrix languages, the product operator `*` operates elementwise in NumPy arrays. The matrix product can be performed using the `dot` function or method:

```
>>>                                                                      >>>
>>> A = np.array( [[1,1],
...                [0,1]] )
>>> B = np.array( [[2,0],
...                [3,4]] )
>>> A*B                              # elementwise product
array([[2, 0],
       [0, 4]])
>>> A.dot(B)                         # matrix product
array([[5, 4],
       [3, 4]])
>>> np.dot(A, B)                     # another matrix product
array([[5, 4],
       [3, 4]])
```

Many unary operations, such as computing the sum of all the elements in the array, are implemented as methods of the `ndarray` class.

```
>>> a = np.random.random((2,3))
>>> a
array([[ 0.18626021,  0.34556073,  0.39676747],
       [ 0.53881673,  0.41919451,  0.6852195 ]])
>>> a.sum()
2.5718191614547998
>>> a.min()
0.1862602113776709
>>> a.max()
0.6852195003967595
```

By default, these operations apply to the array as though it were a list of numbers, regardless of its shape. However, by specifying the `axis` parameter you can apply an operation along the specified axis of an array:

```
>>> b = np.arange(12).reshape(3,4)
>>> b
array([[ 0,  1,  2,  3],
       [ 4,  5,  6,  7],
       [ 8,  9, 10, 11]])
>>>
>>> b.sum(axis=0)                          # sum of each column
array([12, 15, 18, 21])
>>>
>>> b.min(axis=1)                          # min of each row
array([0, 4, 8])
>>>
>>> b.cumsum(axis=1)                       # cumulative sum along each row
array([[ 0,  1,  3,  6],
       [ 4,  9, 15, 22],
       [ 8, 17, 27, 38]])
```

# Indexing, Slicing and Iterating

**One-dimensional** arrays can be indexed, sliced and iterated over, much like lists and other Python sequences.

```
>>> a = np.arange(10)**3
>>> a
array([  0,   1,   8,  27,  64, 125, 216, 343, 512, 729])
>>> a[2]
8
>>> a[2:5]
array([ 8, 27, 64])
>>> a[:6:2] = -1000    # equivalent to a[0:6:2] = -1000; from start to position 6, exclusive, set e
very 2nd element to -1000
>>> a
array([-1000,     1, -1000,    27, -1000,   125,   216,   343,   512,   729])
>>> a[ : :-1]                                 # reversed a
array([  729,   512,   343,   216,   125, -1000,    27, -1000,     1, -1000])
>>> for i in a:
...     print(i**(1/3.))
...
nan
1.0
nan
3.0
nan
5.0
6.0
7.0
8.0
9.0
```

# Changing the shape of an array

An array has a shape given by the number of elements along each axis:

```
>>> a = np.floor(10*np.random.random((3,4)))
>>> a
array([[ 2.,  8.,  0.,  6.],
       [ 4.,  5.,  1.,  1.],
       [ 8.,  9.,  3.,  6.]])
>>> a.shape
(3, 4)
```

The shape of an array can be changed with various commands. Note that the following three commands all return a modified array, but do not change the original array:

```
>>> a.ravel()  # returns the array, flattened
array([ 2.,  8.,  0.,  6.,  4.,  5.,  1.,  1.,  8.,  9.,  3.,  6.])
>>> a.reshape(6,2)  # returns the array with a modified shape
array([[ 2.,  8.],
       [ 0.,  6.],
       [ 4.,  5.],
       [ 1.,  1.],
       [ 8.,  9.],
       [ 3.,  6.]])
>>> a.T  # returns the array, transposed
array([[ 2.,  4.,  8.],
       [ 8.,  5.,  9.],
       [ 0.,  1.,  3.],
       [ 6.,  1.,  6.]])
>>> a.T.shape
(4, 3)
>>> a.shape
(3, 4)
```

# Stacking together different arrays

Several arrays can be stacked together along different axes:

```
>>> a = np.floor(10*np.random.random((2,2)))
>>> a
array([[ 8.,  8.],
       [ 0.,  0.]])
>>> b = np.floor(10*np.random.random((2,2)))
>>> b
array([[ 1.,  8.],
       [ 0.,  4.]])
>>> np.vstack((a,b))
array([[ 8.,  8.],
       [ 0.,  0.],
       [ 1.,  8.],
       [ 0.,  4.]])
>>> np.hstack((a,b))
array([[ 8.,  8.,  1.,  8.],
       [ 0.,  0.,  0.,  4.]])
```

>>>

# Copies and Views

When operating and manipulating arrays, their data is sometimes copied into a new array and sometimes not. This is often a source of confusion for beginners. There are three cases:

## No Copy at All

Simple assignments make no copy of array objects or of their data.

```
>>> a = np.arange(12)
>>> b = a                # no new object is created
>>> b is a               # a and b are two names for the same ndarray object
True
>>> b.shape = 3,4        # changes the shape of a
>>> a.shape
(3, 4)
```

# View or Shallow Copy

Different array objects can share the same data. The `view` method creates a new array object that looks at the same data.

```
>>>
>>> c = a.view()
>>> c is a
False
>>> c.base is a                         # c is a view of the data owned by a
True
>>> c.flags.owndata
False
>>>
>>> c.shape = 2,6                        # a's shape doesn't change
>>> a.shape
(3, 4)
>>> c[0,4] = 1234                        # a's data changes
>>> a
array([[   0,    1,    2,    3],
       [1234,    5,    6,    7],
       [   8,    9,   10,   11]])
```

# Deep Copy

The `copy` method makes a complete copy of the array and its data.

```
>>> d = a.copy()                              # a new array object with new data is created
>>> d is a
False
>>> d.base is a                               # d doesn't share anything with a
False
>>> d[0,0] = 9999
>>> a
array([[   0,   10,   10,    3],
       [1234,   10,   10,    7],
       [   8,   10,   10,   11]])
```

# Fancy indexing and index tricks

NumPy offers more indexing facilities than regular Python sequences. In addition to indexing by integers and slices, as we saw before, arrays can be indexed by arrays of integers and arrays of booleans.

## Indexing with Arrays of Indices

```
                                                                              >>>
>>> a = np.arange(12)**2                    # the first 12 square numbers
>>> i = np.array( [ 1,1,3,8,5 ] )           # an array of indices
>>> a[i]                                     # the elements of a at the positions i
array([ 1,  1,  9, 64, 25])
>>>
>>> j = np.array( [ [ 3, 4], [ 9, 7 ] ] )   # a bidimensional array of indices
>>> a[j]                                     # the same shape as j
array([[ 9, 16],
       [81, 49]])
```

# Scipy

# Introduction

## Contents

- Introduction
  - SciPy Organization
  - Finding Documentation

SciPy is a collection of mathematical algorithms and convenience functions built on the Numpy extension of Python. It adds significant power to the interactive Python session by providing the user with high-level commands and classes for manipulating and visualizing data. With SciPy an interactive Python session becomes a data-processing and system-prototyping environment rivaling systems such as MATLAB, IDL, Octave, R-Lab, and SciLab.

# SciPy Organization

SciPy is organized into subpackages covering different scientific computing domains. These are summarized in the following table:

| Subpackage | Description |
| --- | --- |
| `cluster` | Clustering algorithms |
| `constants` | Physical and mathematical constants |
| `fftpack` | Fast Fourier Transform routines |
| `integrate` | Integration and ordinary differential equation solvers |
| `interpolate` | Interpolation and smoothing splines |
| `io` | Input and Output |
| `linalg` | Linear algebra |
| `ndimage` | N-dimensional image processing |
| `odr` | Orthogonal distance regression |
| `optimize` | Optimization and root-finding routines |
| `signal` | Signal processing |
| `sparse` | Sparse matrices and associated routines |
| `spatial` | Spatial data structures and algorithms |
| `special` | Special functions |
| `stats` | Statistical distributions and functions |

Scipy sub-packages need to be imported separately, for example:

```
>>> from scipy import linalg, optimize
```

# HW 3

- Write a short tutorial on numpy and scipy

- Cover basic functionality with worked examples

- Due Sunday by 11:59 pm


- Evaluations due Tuesday by 11:59pm

  - Use the rubric!

  - Justify your scores on each section of the rubric!

# EXAMPLE!
## (actual notebook)

# In-class problems

- Import the numpy package under the name np

- Create a zeros vector of size 20

- Create a vector with values ranging from 15 to 41

- Reverse a vector (first element becomes last)

- Create a 4x4 matrix with values ranging from 0 to 15

- Find indices of non-zero elements from [1,2,0,0,4,0]

- Create a 3x3x2 array with random values

- Given a 1D array, negate all elements which are between 3 and 8, in place.

- Create random vector of size 10 and replace the maximum value by 0

# HW 3

- Write a short tutorial on numpy and scipy

- Cover basic functionality with worked examples

- Due Sunday by 11:59 pm


- Evaluations due Tuesday by 11:59pm

  - Use the rubric!

  - Justify your scores on each section of the **rubric**!