# Introduction to numpy:

## Package for scientific computing with Python

Numerical Python, or "Numpy" for short, is a foundational package on which many of the most common data science packages are built. Numpy provides us with high performance multi-dimensional arrays which we can use as vectors or matrices.

The key features of numpy are:

- ndarrays: n-dimensional arrays of the same data type which are fast and space-efficient. There are a number of built-in methods for ndarrays which allow for rapid processing of data without using loops (e.g., compute the mean).
- Broadcasting: a useful tool which defines implicit behavior between multi-dimensional arrays of different sizes.
- Vectorization: enables numeric operations on ndarrays.
- Input/Output: simplifies reading and writing of data from/to file.

**Additional Recommended Resources:**
Numpy Documentation (https://docs.scipy.org/doc/numpy/reference/)
*Python for Data Analysis* by Wes McKinney
*Python Data science Handbook* by Jake VanderPlas

# Getting started with ndarray

**ndarrays** are time and space-efficient multidimensional arrays at the core of numpy. Like the data structures in Week 2, let's get started by creating ndarrays using the numpy package.

## How to create Rank 1 numpy arrays:

In [ ]:

```python
import numpy as np

an_array = np.array([3, 33, 333])  # Create a rank 1 array

print(type(an_array))              # The type of an ndarray is: "<class 'numpy.ndarray'>"
```

In [ ]:

```
# test the shape of the array we just created, it should have just one dimension
(Rank 1)
print(an_array.shape)
```

In [ ]:

```
# because this is a 1-rank array, we need only one index to accesss each element
print(an_array[0], an_array[1], an_array[2])
```

In [ ]:

```
an_array[0] =888           # ndarrays are mutable, here we change an element of
the array

print(an_array)
```

# How to create a Rank 2 numpy array:

A rank 2 **ndarray** is one with two dimensions. Notice the format below of [ [row] , [row] ]. 2 dimensional arrays are great for representing matrices which are often useful in data science.

In [ ]:

```
another = np.array([[11,12,13],[21,22,23]])    # Create a rank 2 array

print(another)  # print the array

print("The shape is 2 rows, 3 columns: ", another.shape)  # rows x columns

print("Accessing elements [0,0], [0,1], and [1,0] of the ndarray: ", another[0,
0], ", ",another[0, 1],", ", another[1, 0])
```

# There are many way to create numpy arrays:

Here we create a number of different size arrays with different shapes and different pre-filled values. numpy has a number of built in methods which help us quickly and easily create multidimensional arrays.

In [ ]:

```
import numpy as np

# create a 2x2 array of zeros
ex1 = np.zeros((2,2))
print(ex1)
```

In [ ]:

```
# create a 2x2 array filled with 9.0
ex2 = np.full((2,2), 9.0)
print(ex2)
```

In [ ]:

```python
# create a 2x2 matrix with the diagonal 1s and the others 0
ex3 = np.eye(2,2)
print(ex3)
```

In [ ]:

```python
# create an array of ones
ex4 = np.ones((1,2))
print(ex4)
```

In [ ]:

```python
# notice that the above ndarray (ex4) is actually rank 2, it is a 2x1 array
print(ex4.shape)

# which means we need to use two indexes to access an element
print()
print(ex4[0,1])
```

In [ ]:

```python
# create an array of random floats between 0 and 1
ex5 = np.random.random((2,2))
print(ex5)
```

# Array Indexing

## Slice indexing:

Similar to the use of slice indexing with lists and strings, we can use slice indexing to pull out sub-regions of ndarrays.

In [ ]:

```python
import numpy as np

# Rank 2 array of shape (3, 4)
an_array = np.array([[11,12,13,14], [21,22,23,24], [31,32,33,34]])
print(an_array)
```

Use array slicing to get a subarray consisting of the first 2 rows x 2 columns.

In [ ]:

```python
a_slice = an_array[:2, 1:3]
print(a_slice)
```

When you modify a slice, you actually modify the underlying array.

In [ ]:

```python
print("Before:", an_array[0, 1])    #inspect the element at 0, 1
a_slice[0, 0] = 1000    # a_slice[0, 0] is the same piece of data as an_array[0, 1]
print("After:", an_array[0, 1])
```

# Use both integer indexing & slice indexing

We can use combinations of integer indexing and slice indexing to create different shaped matrices.

In [ ]:

```python
# Create a Rank 2 array of shape (3, 4)
an_array = np.array([[11,12,13,14], [21,22,23,24], [31,32,33,34]])
print(an_array)
```

In [ ]:

```python
# Using both integer indexing & slicing generates an array of lower rank
row_rank1 = an_array[1, :]    # Rank 1 view

print(row_rank1, row_rank1.shape)  # notice only a single []
```

In [ ]:

```python
# Slicing alone: generates an array of the same rank as the an_array
row_rank2 = an_array[1:2, :]  # Rank 2 view

print(row_rank2, row_rank2.shape)   # Notice the [[ ]]
```

In [ ]:

```python
#We can do the same thing for columns of an array:

print()
col_rank1 = an_array[:, 1]
col_rank2 = an_array[:, 1:2]

print(col_rank1, col_rank1.shape)  # Rank 1
print()
print(col_rank2, col_rank2.shape)  # Rank 2
```

# Array Indexing for changing elements:

Sometimes it's useful to use an array of indexes to access or change elements.

In [ ]:

```python
# Create a new array
an_array = np.array([[11,12,13], [21,22,23], [31,32,33], [41,42,43]])

print('Original Array:')
print(an_array)
```

In [ ]:

```python
# Create an array of indices
col_indices = np.array([0, 1, 2, 0])
print('\nCol indices picked : ', col_indices)

row_indices = np.arange(4)
print('\nRows indices picked : ', row_indices)
```

In [ ]:

```python
# Examine the pairings of row_indices and col_indices.  These are the elements we'll change next.
for row,col in zip(row_indices,col_indices):
    print(row, ", ",col)
```

In [ ]:

```python
# Select one element from each row
print('Values in the array at those indices: ',an_array[row_indices, col_indices])
```

In [ ]:

```python
# Change one element from each row using the indices selected
an_array[row_indices, col_indices] += 100000

print('\nChanged Array:')
print(an_array)
```

# Boolean Indexing

## Array Indexing for changing elements:

In [ ]:

```python
# create a 3x2 array
an_array = np.array([[11,12], [21, 22], [31, 32]])
print(an_array)
```

In [ ]:

```python
# create a filter which will be boolean values for whether each element meets this condition
filter = (an_array > 15)
filter
```

Notice that the filter is a same size ndarray as an_array which is filled with True for each element whose corresponding element in an_array which is greater than 15 and False for those elements whose value is less than 15.

In [ ]:

```python
# we can now select just those elements which meet that criteria
print(an_array[filter])
```

In [ ]:

```python
# For short, we could have just used the approach below without the need for the
separate filter array.

an_array[an_array > 15]
```

What is particularly useful is that we can actually change elements in the array applying a similar logical filter. Let's add 100 to all the even values.

In [ ]:

```python
an_array[an_array % 2 == 0] +=100
print(an_array)
```

# Datatypes and Array Operations

## Datatypes:

In [ ]:

```python
ex1 = np.array([11, 12]) # Python assigns the  data type
print(ex1.dtype)
```

In [ ]:

```python
ex2 = np.array([11.0, 12.0]) # Python assigns the  data type
print(ex2.dtype)
```

In [ ]:

```python
ex3 = np.array([11, 21], dtype=np.int64) #You can also tell Python the  data typ
e
print(ex3.dtype)
```

In [ ]:

```python
# you can use this to force floats into integers (using floor function)
ex4 = np.array([11.1,12.7], dtype=np.int64)
print(ex4.dtype)
print()
print(ex4)
```

In [ ]:

```python
# you can use this to force integers into floats if you anticipate
# the values may change to floats later
ex5 = np.array([11, 21], dtype=np.float64)
print(ex5.dtype)
print()
print(ex5)
```

# Arithmetic Array Operations:

In [ ]:

```python
x = np.array([[111,112],[121,122]], dtype=np.int)
y = np.array([[211.1,212.1],[221.1,222.1]], dtype=np.float64)

print(x)
print()
print(y)
```

In [ ]:

```python
# add
print(x + y)              # The plus sign works
print()
print(np.add(x, y))   # so does the numpy function "add"
```

In [ ]:

```python
# subtract
print(x - y)
print()
print(np.subtract(x, y))
```

In [ ]:

```python
# multiply
print(x * y)
print()
print(np.multiply(x, y))
```

In [ ]:

```python
# divide
print(x / y)
print()
print(np.divide(x, y))
```

In [ ]:

```python
# square root
print(np.sqrt(x))
```

In [ ]:

```python
# exponent (e ** x)
print(np.exp(x))
```

# Statistical Methods, Sorting, and Set Operations:

## Basic Statistical Operations:

In [ ]:

```python
# setup a random 2 x 4 matrix
arr = 10 * np.random.randn(2,5)
print(arr)
```

In [ ]:

```python
# compute the mean for all elements
print(arr.mean())
```

In [ ]:

```python
# compute the means by row
print(arr.mean(axis = 1))
```

In [ ]:

```python
# compute the means by column
print(arr.mean(axis = 0))
```

In [ ]:

```python
# sum all the elements
print(arr.sum())
```

In [ ]:

```python
# compute the medians
print(np.median(arr, axis = 1))
```

## Sorting:

In [ ]:

```python
# create a 10 element array of randoms
unsorted = np.random.randn(10)

print(unsorted)
```

In [ ]:

```python
# create copy and sort
sorted = np.array(unsorted)
sorted.sort()

print(sorted)
print()
print(unsorted)
```

In [ ]:

```python
# inplace sorting
unsorted.sort()

print(unsorted)
```

## Finding Unique elements:

In [ ]:

```python
array = np.array([1,2,1,4,2,1,4,2])

print(np.unique(array))
```

## Set Operations with np.array data type:

In [ ]:

```python
s1 = np.array(['desk','chair','bulb'])
s2 = np.array(['lamp','bulb','chair'])
print(s1, s2)
```

In [ ]:

```python
print( np.intersect1d(s1, s2) )
```

In [ ]:

```python
print( np.union1d(s1, s2) )
```

In [ ]:

```python
print( np.setdiff1d(s1, s2) )# elements in s1 that are not in s2
```

In [ ]:

```python
print( np.in1d(s1, s2) )#which element of s1 is also in s2
```

# Broadcasting:

Introduction to broadcasting.
For more details, please see:
https://docs.scipy.org/doc/numpy-1.10.1/user/basics.broadcasting.html (https://docs.scipy.org/doc/numpy-1.10.1/user/basics.broadcasting.html)

In [ ]:

```python
import numpy as np

start = np.zeros((4,3))
print(start)
```

In [ ]:

```python
# create a rank 1 ndarray with 3 values
add_rows = np.array([1, 0, 2])
print(add_rows)
```

In [ ]:

```python
y = start + add_rows    # add to each row of 'start' using broadcasting
print(y)
```

In [ ]:

```python
# create an ndarray which is 4 x 1 to broadcast across columns
add_cols = np.array([[0,1,2,3]])
add_cols = add_cols.T

print(add_cols)
```

In [ ]:

```python
# add to each column of 'start' using broadcasting
y = start + add_cols
print(y)
```

In [ ]:

```python
# this will just broadcast in both dimensions
add_scalar = np.array([1])
print(start+add_scalar)
```

Example from the slides:

In [ ]:

```python
# create our 3x4 matrix
arrA = np.array([[1,2,3,4],[5,6,7,8],[9,10,11,12]])
print(arrA)
```

In [ ]:

```python
# create our 4x1 array
arrB = [0,1,0,2]
print(arrB)
```

In [ ]:

```python
# add the two together using broadcasting
print(arrA + arrB)
```

# Speedtest: ndarrays vs lists

First setup paramaters for the speed test. We'll be testing time to sum elements in an ndarray versus a list.

In [ ]:

```python
from numpy import arange
from timeit import Timer

size    = 1000000
timeits = 1000
```

In [ ]:

```python
# create the ndarray with values 0,1,2...,size-1
nd_array = arange(size)
print( type(nd_array) )
```

In [ ]:

```python
# timer expects the operation as a parameter,
# here we pass nd_array.sum()
timer_numpy = Timer("nd_array.sum()", "from __main__ import nd_array")

print("Time taken by numpy ndarray: %f seconds" %
      (timer_numpy.timeit(timeits)/timeits))
```

In [ ]:

```python
# create the list with values 0,1,2...,size-1
a_list = list(range(size))
print (type(a_list) )
```

In [ ]:

```python
# timer expects the operation as a parameter, here we pass sum(a_list)
timer_list = Timer("sum(a_list)", "from __main__ import a_list")

print("Time taken by list:  %f seconds" %
      (timer_list.timeit(timeits)/timeits))
```

# Read or Write to Disk:

## Binary Format:

In [ ]:

```python
x = np.array([ 23.23, 24.24] )
```

In [ ]:

```python
np.save('an_array', x)
```

In [ ]:

```python
np.load('an_array.npy')
```

## Text Format:

In [ ]:

```python
np.savetxt('array.txt', X=x, delimiter=',')
```

In [ ]:

```python
## For Windows replace "cat" with "type"
!cat array.txt
```

In [ ]:

```python
np.loadtxt('array.txt', delimiter=',')
```

# Additional Common ndarray Operations

## Dot Product on Matrices and Inner Product on Vectors:

In [ ]:

```python
# determine the dot product of two matrices
x2d = np.array([[1,1],[1,1]])
y2d = np.array([[2,2],[2,2]])

print(x2d.dot(y2d))
print()
print(np.dot(x2d, y2d))
```

In [ ]:

```python
# determine the inner product of two vectors
a1d = np.array([9 , 9 ])
b1d = np.array([10, 10])

print(a1d.dot(b1d))
print()
print(np.dot(a1d, b1d))
```

In [ ]:

```python
# dot produce on an array and vector
print(x2d.dot(a1d))
print()
print(np.dot(x2d, a1d))
```

## Sum:

In [ ]:

```python
# sum elements in the array
ex1 = np.array([[11,12],[21,22]])

print(np.sum(ex1))          # add all members
```

In [ ]:

```python
print(np.sum(ex1, axis=0))  # columnwise sum
```

In [ ]:

```python
print(np.sum(ex1, axis=1))  # rowwise sum
```

## Element-wise Functions:

For example, let's compare two arrays values to get the maximum of each.

In [ ]:

```python
# random array
x = np.random.randn(8)
x
```

In [ ]:

```python
# another random array
y = np.random.randn(8)
y
```

In [ ]:

```python
# returns element wise maximum between two arrays

np.maximum(x, y)
```

## Reshaping array:

In [ ]:

```python
# grab values from 0 through 19 in an array
arr = np.arange(20)
print(arr)
```

In [ ]:

```python
# reshape to be a 4 x 5 matrix
arr.reshape(4,5)
```

## Transpose:

In [ ]:

```python
# transpose
ex1 = np.array([[11,12],[21,22]])

ex1.T
```

## Indexing using where():

In [ ]:

```python
x_1 = np.array([1,2,3,4,5])

y_1 = np.array([11,22,33,44,55])

filter = np.array([True, False, True, False, True])
```

In [ ]:

```python
out = np.where(filter, x_1, y_1)
print(out)
```

In [ ]:

```python
mat = np.random.rand(5,5)
mat
```

In [ ]:

```python
np.where( mat > 0.5, 1000, -1)
```

## "any" or "all" conditionals:

In [ ]:

```python
arr_bools = np.array([ True, False, True, True, False ])
```

In [ ]:

```python
arr_bools.any()
```

In [ ]:

```python
arr_bools.all()
```

# Random Number Generation:

In [ ]:

```python
Y = np.random.normal(size = (1,5))[0]
print(Y)
```

In [ ]:

```python
Z = np.random.randint(low=2,high=50,size=4)
print(Z)
```

In [ ]:

```python
np.random.permutation(Z) #return a new ordering of elements in Z
```

In [ ]:

```python
np.random.uniform(size=4) #uniform distribution
```

In [ ]:

```python
np.random.normal(size=4) #normal distribution
```

# Merging data sets:

In [ ]:

```python
K = np.random.randint(low=2,high=50,size=(2,2))
print(K)

print()
M = np.random.randint(low=2,high=50,size=(2,2))
print(M)
```

In [ ]:

```python
np.vstack((K,M))
```

In [ ]:

```python
np.hstack((K,M))
```

In [ ]:

```python
np.concatenate([K, M], axis = 0)
```

In [ ]:

```python
np.concatenate([K, M.T], axis = 1)
```

In [ ]: