

# Fundamentals of Information Systems

## Python Programming (for Data Science)

### Master's Degree in Data Science

Gabriele Tolomei

[gtolomei@math.unipd.it](mailto:gtolomei@math.unipd.it)

University of Padua, Italy

2018/2019

November, 8 2018

# Lecture 6: Numerical Python (`numpy`)

# What is **numpy**?

- It stands for **num**erical **py**thon, and is one of the core packages for numerical/scientific computing in Python.
- Most computational packages providing scientific functionality use **numpy array objects** as the building block for data exchange.
- You can find more about **numpy** on the official [website](#).

In [1]:

```
1 """
2 As any other third-party module, the numpy module has to be imported before it can be used.
3 If you installed Python with Anaconda, numpy would be just available to you.
4 This is usually how numpy is imported and aliased. Although you could also
5 use another syntax like 'from numpy import *', I strongly encourage you to define an alias,
6 as this will help you to identify numpy's functions in your code.
7 """
8 import numpy as np
```

# What is inside `numpy`?

- `ndarray`: an efficient multi-dimensional array providing fast array-oriented arithmetic operations.
- Mathematical functions for fast operations on entire arrays of data without having to write loops.
- Tools for reading array data from (writing array data to) disk and working with memory-mapped files.
- Linear algebra, random number generation, and Fourier transform capabilities.
- A C API for connecting `numpy` with libraries written in C, C++, or FORTRAN.

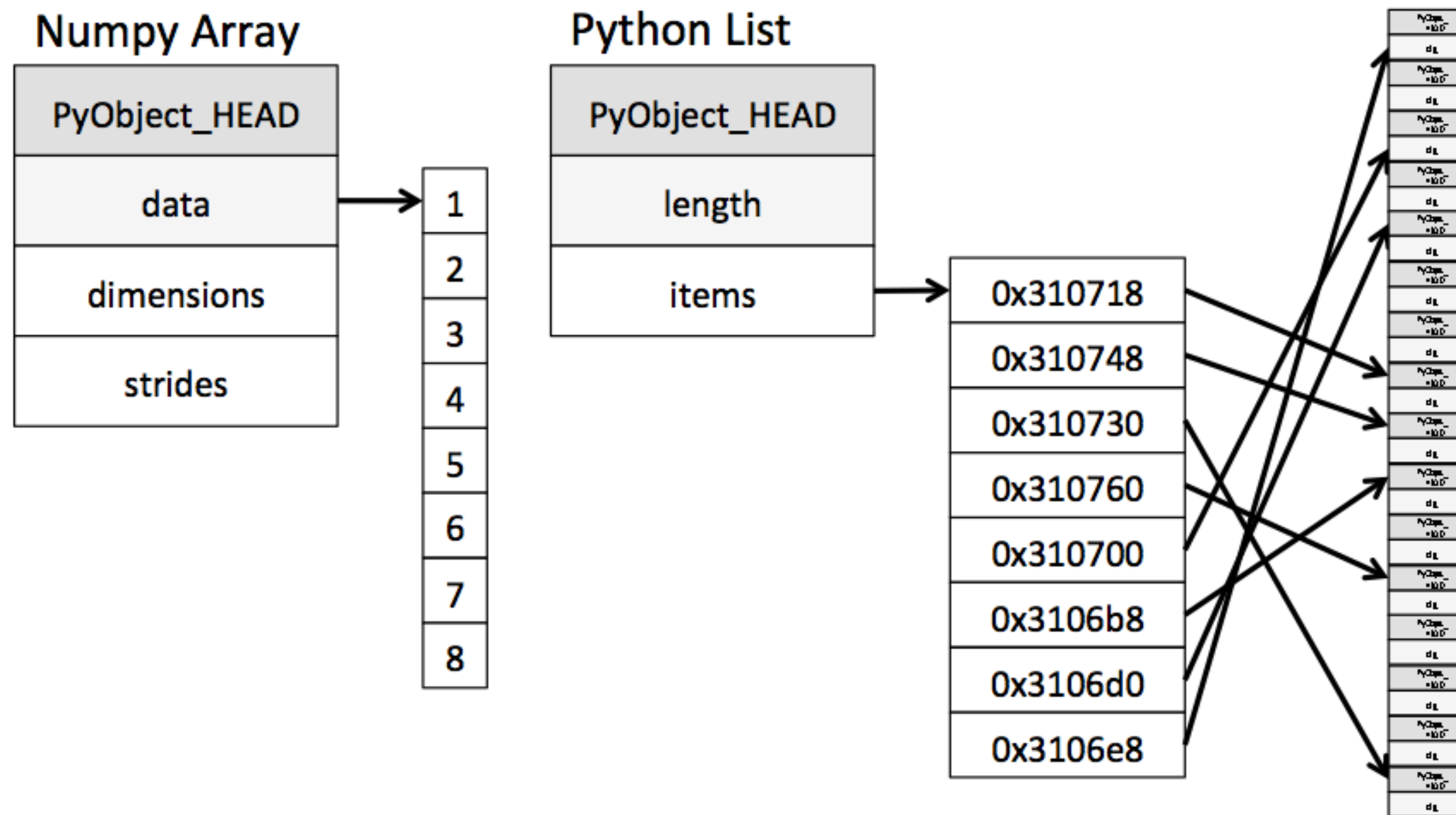
# General purpose `numpy`

- Because `numpy` provides an easy-to-use C API, it is straightforward to pass data to/from external libraries written in a low-level language.
- `numpy` itself does not provide modeling nor scientific functionality, but knowing of `numpy` basics will help you use tools with array-oriented semantics, like `pandas`.
- In this class, we will in fact use `pandas`, which is tailored to tabular data and also provides some more domain-specific functionality like time series manipulation, which is not present in `numpy`.

# Space Efficiency of `numpy`'s `ndarray`

- `numpy`'s importance for numerical computations in Python is due to its design for efficiency (especially when operating on large arrays of data).
- It internally stores data in a **contiguous** block of memory, independent of other built-in Python objects.

# Space Efficiency of `numpy`'s `ndarray`





# Time Efficiency of `numpy`'s `ndarray`

- `ndarray`'s efficient memory occupation implies also computational time efficiency.
- Its library of algorithms mostly written in low-level C can operate on this memory without introducing any overhead due to type checking.
- `numpy` operations perform complex computations on entire arrays without the need for Python `for` loops (i.e., knowing the address of the memory block and the data type, it is just simple arithmetic).
- Spatial locality in memory access patterns results in performance gains notably due to the CPU cache (sequential locality, or locality of reference).
- Since items are stored contiguously in memory, `numpy` can take advantage of **vectorized instructions** provided by modern CPUs.

# Efficiency of `numpy`: a real example

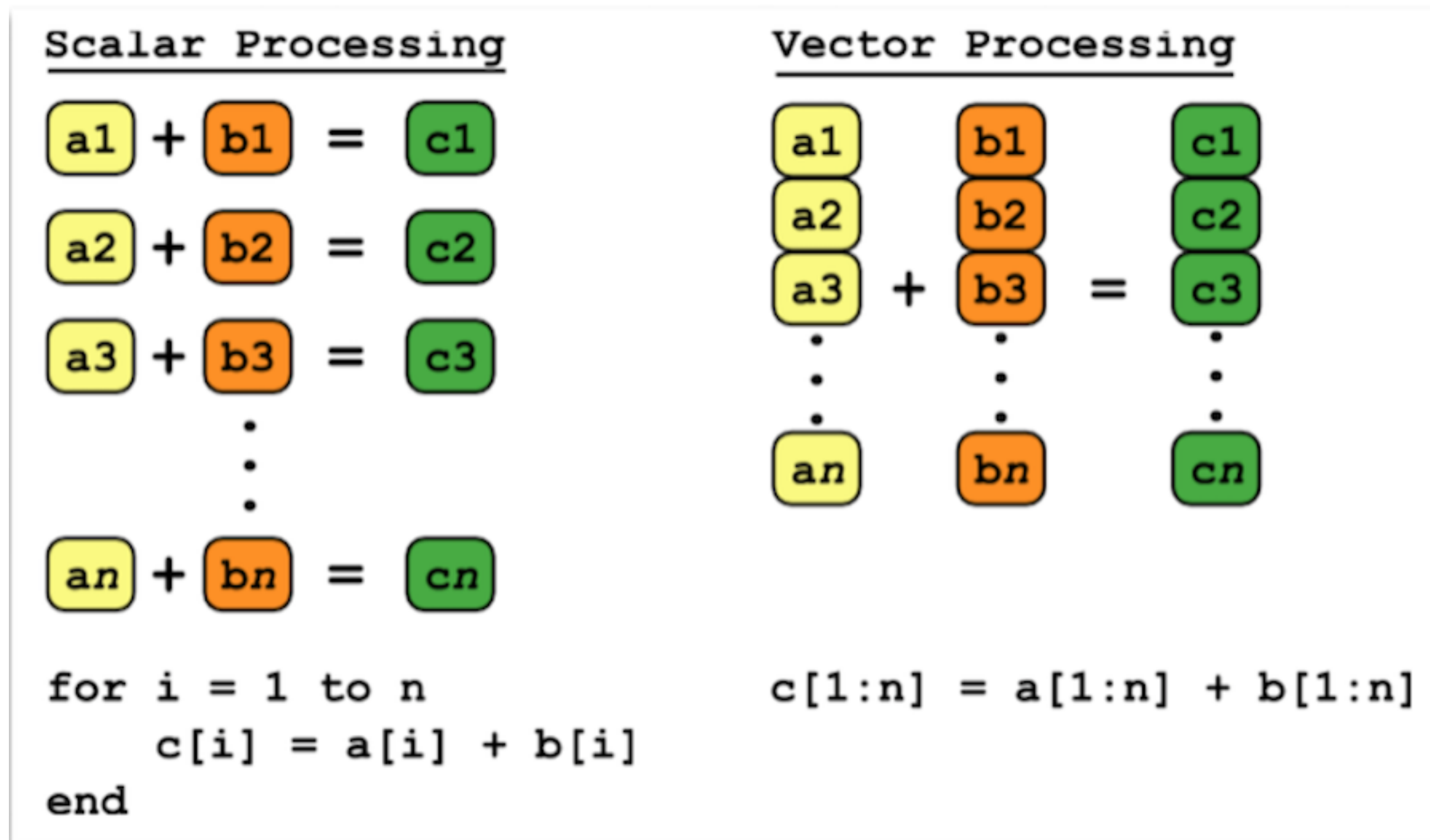
- To validate the efficiency of `numpy` in contrast with built-in Python list, just try to run the code snippet below:

```
# create a numpy array with 1M integers
my_arr = np.arange(1000000)
# create a built-in list with 1M integers
my_list = list(range(1000000))
# double each element of the numpy array
my_arr2 = my_arr * 2
# double each element of the built-in list
my_list2 = [x * 2 for x in my_list]
```

- `numpy`-based algorithms are expected to be **10 to 100** times faster than their pure Python counterparts and use significantly less memory.

# Scalar vs. Vector Processing

Vector processing is also known as **Single Instruction Multiple Data (SIMD)**



**ndarray: A Multidimensional Array Object**

# Properties of `ndarray`

- A fast, flexible, generic multidimensional container for large homogeneous data sets in Python.
- Arrays enable you to perform mathematical operations on whole blocks of data using similar syntax to the equivalent operations between scalars.
- All the elements of an `ndarray` must be of the **same** type.
- Every array has a `shape`, a tuple indicating the size of each dimension, and a `dtype`, an object describing the data type of the array.

In [2]:

```
# Generate some random data over a 2x3 array (i.e., a matrix)
data = np.random.randn(2, 3)
# Print out data
print("Original data matrix =\n{}".format(data))
# Multiply each element of the matrix by a constant 10
data10 = data * 10
# Print out the new data matrix
print("data matrix * 10 =\n{}".format(data10))
# Sum two data matrices (i.e., the same as multiply each element by 2)
data2 = data + data
# Print out the new data matrix
print("(data matrix + data matrix) =\n{}".format(data2))
```

Original data matrix =

```
[[ 0.60412659  1.20847991 -1.03453903]
 [ 0.76247938 -0.6180134   1.12874354]]
```

data matrix \* 10 =

```
[[ 6.04126594 12.0847991 -10.34539031]
 [ 7.62479376 -6.18013404 11.28743543]]
```

(data matrix + data matrix) =

```
[[ 1.20825319  2.41695982 -2.06907806]
 [ 1.52495875 -1.23602681  2.25748709]]
```

In [3]:

```
# Showing the shape of the ndarray object  
print("The shape of data is: {}".format(data.shape))  
# Showing the type of objects contained in the ndarray object  
print("The type of objects contained in data is: {}".format(data.dtype))
```

The shape of data is: (2, 3)

The type of objects contained in data is: float64

Creating `ndarray`



In [4]:

```
# Start from a built-in Python list  
data = [42, 2.5, 73, 0, 3, 1.0]  
# The corresponding numpy array can be obtained by calling the np.array function  
arr = np.array(data)  
arr
```

Out[4]: array([ 42. , 2.5, 73. , 0. , 3. , 1. ])

In [5]:

```
1# Nested sequences, like a list of equal-length lists,  
2# will be converted into a multidimensional array  
multi_data = [[1, 2, 3, 4], [5, 6, 7, 8]]  
4# Convert the list of list into a (multidimensional) numpy array  
multi_arr = np.array(multi_data)  
6print("Multidimensional array:\n{}".format(multi_arr))  
7print("Number of dimensions of the array: {}".format(multi_arr.ndim))  
8print("Shape of the array: {}".format(multi_arr.shape))  
9# Unless explicitly specified (more on this later), np.array tries to infer  
10# a good data type for the array that it creates.  
11# The data type is stored in a special dtype metadata object's field.  
12print("Shape of the unidimensional array: {}".format(arr.dtype))  
13print("Shape of the multidimensional array: {}".format(multi_arr.dtype))
```

Multidimensional array:

```
[[1 2 3 4]  
 [5 6 7 8]]
```

Number of dimensions of the array: 2

Shape of the array: (2, 4)

Shape of the unidimensional array: float64

Shape of the multidimensional array: int64

In [6]:

```
1"""
2In addition to np.array, there are a number of other functions for creating new arrays.
3As examples, 'zeros' and 'ones' create arrays of 0's or 1's, respectively,
4with a given length or shape.
5empty' creates an array without initializing its values to any particular value.
6To create a higher dimensional array with these methods, pass a tuple for the shape.
7"""
8print("Creating a unidimensional array with 5 zeros: {}".format(np.zeros(5)))
9print("Creating a multidimensional array (i.e., 3x2 matrix) with all zeros:\n{}\n"
10      .format(np.zeros((3,2))))
11print("Creating two empty multidimensional arrays (i.e., 3x4 matrix):\n{}\n"
12      .format(np.empty((2, 3, 4))))
```

Creating a unidimensional array with 5 zeros: [ 0. 0. 0. 0. 0.]

Creating a multidimensional array (i.e., 3x2 matrix) with all zeros:

```
[[ 0.  0.]
 [ 0.  0.]
 [ 0.  0.]]
```

Creating two empty multidimensional arrays (i.e., 3x4 matrix):

```
[[[ 2.31584178e+077  2.31584178e+077  6.42285340e-323  0.00000000e+000]
 [ 0.00000000e+000  0.00000000e+000  0.00000000e+000  0.00000000e+000]
 [ 0.00000000e+000  0.00000000e+000  0.00000000e+000  0.00000000e+000]]

 [[ 0.00000000e+000  0.00000000e+000  0.00000000e+000  1.14411728e-308]
 [ 2.31584178e+077  2.31584178e+077  2.47032823e-323  0.00000000e+000]
 [ 0.00000000e+000  0.00000000e+000  2.31584178e+077  2.00390288e+000]]]
```

```
In [7]: 1 # 'arange' is an array-valued version of the built-in Python 'range' function  
        2 np.arange(16)
```

```
Out[7]: array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15])
```

# Table of `numpy` Functions to Create `ndarray`

Function	Description
<code>array</code>	Convert input data (list, tuple, array, or other sequence type) to an ndarray either by inferring a dtype or explicitly specifying a dtype. Copies the input data by default.
<code>asarray</code>	Convert input to ndarray, but do not copy if the input is already an ndarray
<code>arange</code>	Like the built-in <code>range</code> but returns an ndarray instead of a list.
<code>ones</code> , <code>ones_like</code>	Produce an array of all 1's with the given shape and dtype. <code>ones_like</code> takes another array and produces a ones array of the same shape and dtype.
<code>zeros</code> , <code>zeros_like</code>	Like <code>ones</code> and <code>ones_like</code> but producing arrays of 0's instead
<code>empty</code> , <code>empty_like</code>	Create new arrays by allocating new memory, but do not populate with any values like <code>ones</code> and <code>zeros</code>
<code>full</code> , <code>full_like</code>	Produce an array of the given shape and dtype with all values set to the indicated "fill value". <code>full_like</code> takes another array and produces a filled array of the same shape and dtype.
<code>eye</code> , <code>identity</code>	Create a square N x N identity matrix (1's on the diagonal and 0's elsewhere)

# Data Types for `ndarray`

# Data Type: `dtype`

- The data type or `dtype` is a special object containing the information (or `metadata`) that the ndarray needs to interpret a chunk of memory as a particular type of data.
- In most cases, `dtype`s provide a mapping directly onto an underlying disk or memory representation, which makes it easy to read and write binary streams of data.
- Numerical dtypes are named the same way as built-in numerics, yet they also contain the number of bits per element. E.g., `float64` is the `numpy` equivalent of a standard double-precision floating point.

In [8]:

```
1 # Explicitly declare the dtype of the array at definition time  
2 # float64  
arr1 = np.array([1, 2, 3], dtype=np.float64)  
4 # int32  
arr2 = np.array([1, 2, 3], dtype=np.int32)  
print("arr1 data type is: {}".format(arr1.dtype))  
print("arr2 data type is: {}".format(arr2.dtype))
```

arr1 data type is: float64

arr2 data type is: int32



# Table of dtype (1 of 2)

Type	Type Code	Description
<code>int8, uint8</code>	<code>i1, u1</code>	Signed and unsigned 8-bit (1 byte) integer types
<code>int16, uint16</code>	<code>i2, u2</code>	Signed and unsigned 16-bit integer types
<code>int32, uint32</code>	<code>i4, u4</code>	Signed and unsigned 32-bit integer types
<code>int64, uint64</code>	<code>i8, u8</code>	Signed and unsigned 32-bit integer types
<code>float16</code>	<code>f2</code>	Half-precision floating point
<code>float32</code>	<code>f4</code> or <code>f</code>	Standard single-precision floating point. Compatible with C <code>float</code>
<code>float64</code>	<code>f8</code> or <code>d</code>	Standard double-precision floating point. Compatible with C <code>double</code> and Python <code>float</code> object
<code>float128</code>	<code>f16</code> or <code>g</code>	Extended-precision floating point

# Table of dtype (2 of 2)

<code>complex64,</code> <code>complex128,</code> <code>complex256</code>	<code>c8,</code> <code>c16,</code> <code>c32</code>	Complex numbers represented by two 32, 64, or 128 floats, respectively
<code>bool</code>	<code>?</code>	Boolean type storing <code>True</code> and <code>False</code> values
<code>object</code>	<code>O</code>	Python object type, a value can be any Python object
<code>string_</code>	<code>S</code>	Fixed-length ASCII string type (1 byte per character). For example, to create a string dtype with length 10, use <code>'S10'</code> .
<code>unicode_</code>	<code>U</code>	Fixed-length unicode type (number of bytes platform specific). Same specification semantics as <code>string_</code> (e.g. <code>'U10'</code> ).

Casting to a specific `dtype` using `astype`

In [9]:

```
1"""
2Sometimes it may be useful to explicitly convert or cast an array
3from one dtype to another using ndarray's 'astype' method.
4"""
5# Let's define an array using the numpy's array method
6arr = np.array([1, 2, 3, 4, 5])
7print("The (inferred) dtype for the just defined numpy array is: {}".format(arr.dtype))
8# Now, let's convert the inferred dtype (int64) into float64 using 'astype'
9float_arr = arr.astype(np.float64)
10print("The dtype for the cast numpy array is: {}".format(float_arr.dtype))
```

The (inferred) dtype for the just defined numpy array is: int64

The dtype for the cast numpy array is: float64

In [10]:

```
1 """
2 Sometimes it may be useful to explicitly convert or cast an array
3 from one dtype to another using ndarray's 'astype' method.
4 """
5
6 # In the example above, integers are cast to floating point.
7 # What if we cast some floating point numbers to be of integer dtype?
8 # Let's create the numpy array from a list of float numbers
9 arr = np.array([3.7, -1.2, -2.6, 0.5, 12.9, 10.1])
10 print("The original array is: {}".format(arr))
11 print("The array cast to integer is: {}".format(arr.astype(np.int32)))
```

The original array is: [ 3.7 -1.2 -2.6 0.5 12.9 10.1]

The array cast to integer is: [ 3 -1 -2 0 12 10]

In [11]:

```
1"""
2If you have an array of strings representing numbers,
3you can use 'astype' to convert them to numeric form.
4"""
5numeric_strings = np.array(['1.25', '-9.6', '42'], dtype=np.string_)
6# Note that we use 'float' instead of 'np.float64',
7# as numpy aliases the Python types to its own equivalent dtypes.
8print("The original array cast to string is: {}".format(numeric_strings.astype(float)))
9# If casting fails for some reason (like a string that cannot be converted to float64),
10# a ValueError will be raised.
11wrong_numeric_strings = np.array(['1.25', '-9.6', 'h7-25', '42'], dtype=np.string_)
12print("The original array cast to string is: {}".format(wrong_numeric_strings.astype(float)))
```

The original array cast to string is: [ 1.25 -9.6 42. ]

```
-----
ValueError                                Traceback (most recent call last)
<ipython-input-11-e39493940eea> in <module>()
      10 # a ValueError will be raised.
      11 wrong_numeric_strings = np.array(['1.25', '-9.6', 'h7-25', '42'], dtype=np.string_)
--> 12 print("The original array cast to string is: {}".format(wrong_numeric_strings.as-
ype(float)))
```

ValueError: could not convert string to float: 'h7-25'

## About `astype`

Calling `astype` always returns a copy of the original numpy array (even if we apply a "dummy" casting, i.e., if the new `dtype` we want to cast the array to is the same of the original, old `dtype`).

# Operations between arrays and scalars

- `numpy` arrays enables you to express many kinds of "batched" data processing tasks as concise array expressions, instead of writing `for` loops.
- This practice is commonly referred to as **vectorization**.
- In general, vectorized array operations is one or two (or more) orders of magnitude faster than their pure Python equivalents.
- Any arithmetic operations between equal-size arrays applies the operation elementwise.
- Operations between differently sized arrays is called **broadcasting** but won't be further discussed here.



In [12]:

```
# Let's define a simple 2x3 numpy array
arr = np.array([[1., 2., 3.], [4., 5., 6.]])
print("Consider the following {}x{} array:\n{}".\
      .format(arr.shape[0], arr.shape[1], arr))
# Square the values contained in the original array
arr_squared = arr * arr
print("Square the elements of the original array:\n{}".\
      .format(arr_squared))

# Arithmetic operations with scalars are as you would expect,
# propagating the value to each element
reciprocal_arr = 1/arr
print("Compute the reciprocal of the elements of the original array:\n{}".\
      .format(reciprocal_arr))
sqrt_arr = arr ** 0.5
print("Compute the square root of the elements of the original array:\n{}".\
      .format(sqrt_arr))
```

Consider the following 2x3 array:

```
[[ 1.  2.  3.]
 [ 4.  5.  6.]]
```

Square the elements of the original array:

```
[[ 1.  4.  9.]
 [16. 25. 36.]]
```

Compute the reciprocal of the elements of the original array:

```
[[ 1.          0.5          0.33333333]
 [ 0.25        0.2          0.16666667]]
```

Compute the square root of the elements of the original array:

```
[[ 1.          1.41421356  1.73205081]
 [ 2.          2.23606798  2.44948974]]
```

# Basic Indexing and Slicing

In [13]:

```
1 """
2 numpy array indexing is a rich topic, as there are many ways
3 you may want to select a subset of your data or individual elements.
4 One-dimensional arrays are simple; on the surface they act similarly to Python lists.
5 """
6 # Create an ndarray of 10 random integers in the range [0, 50)
7 arr = np.random.randint(low=0, high=50, size=10)
8 print("Original numpy array is: {}".format(arr))
9 print("Accessing the 6-th element of the array: {}".format(arr[5]))
10 print("Extracting from the 6-th to the 8-th element of the array: {}".format(arr[5:8]))
11 # Assigning a new value to a slice is "broadcasted" to all the elements of the slice.
12 arr[5:8] = 12
13 print("Now the numpy array is: {}".format(arr))
```

Original numpy array is: [33 16 49 21 11 17 36 45 35 43]

Accessing the 6-th element of the array: 17

Extracting from the 6-th to the 8-th element of the array: [17 36 45]

Now the numpy array is: [33 16 49 21 11 12 12 12 35 43]

In [14]:

```
1"""
2Differently from Python's built-in lists, numpy array slices are views on the original array.
3This means that the data is not (shallow-)copied,
4and any modifications to the view will be reflected in the source array.
5"""
6
7# Define a Python standard list containing the first 10 non-negative integers [0, 1, ..., 9]
8py_list = [x for x in range(10)]
9
10# Define a numpy array with the same elements
11arr = np.arange(10)
12
13# Slicing the Python list
14sliced_list = py_list[5:8]
15print("Sliced Python list = {}".format(sliced_list))
16
17# Slicing the numpy array
18sliced_arr = arr[5:8]
19print("Sliced numpy array = {}".format(sliced_arr))
20
21# Changing references of the Python sliced list won't change the original list
22sliced_list = [12 for i in sliced_list]
23print("Sliced Python list = {}".format(sliced_list))
24print("Original Python list = {}".format(py_list))
25
26# Changing references of the sliced numpy array will reflect to the original array
27sliced_arr[:] = 12
28print("Sliced numpy array = {}".format(sliced_arr))
29print("Original numpy arr = {}".format(arr))
```

```
Sliced Python list = [5, 6, 7]
Sliced numpy array = [5 6 7]
Sliced Python list = [12, 12, 12]
Original Python list = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
Sliced numpy array = [12 12 12]
Original numpy arr = [ 0  1  2  3  4 12 12 12  8  9]
```

In [15]:

```
1"""
2If you want a copy of a slice of a numpy array instead of a view,
3you will need to explicitly copy the array; for example arr[5:8].copy()
4"""
5arr = np.arange(10)
6sliced_arr = arr[5:8].copy()
7# Changing references of the sliced numpy array will NOT reflect to the original array
8sliced_arr[:] = 12
9print("Sliced numpy array = {}".format(sliced_arr))
10print("Original numpy arr = {}".format(arr))
```

Sliced numpy array = [12 12 12]

Original numpy arr = [0 1 2 3 4 5 6 7 8 9]

In [16]:

```
1 """
2 With higher dimensional arrays, you have many more options.
3 In a two-dimensional array, the elements at each index are no longer scalars
4 but rather one-dimensional arrays.
5 """
6 # Consider the following 3x3 matrix defined as a two-dimensional array
7 arr2d = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
8 print("Original numpy array is:\n{}".format(arr2d))
9 # Accessing the 2-nd element of the matrix above
10 print("The third element of the original array is: {}".format(arr2d[2]))
11 # Thus, individual elements can be accessed recursively.
12 # This is a bit too cumbersome, so you can pass a comma-separated list of indices.
13 print("The third element of the first array is: {}".format(arr2d[0][2]))
14 arr2d[0][2] == arr2d[0, 2]
```

Original numpy array is:

[[1 2 3]

[4 5 6]

[7 8 9]]

The third element of the original array is: [7 8 9]

The third element of the first array is: 3

Out[16]: True

# Indexing on a 2-d Array

		axis 1		
		0	1	2
axis 0	0	0, 0	0, 1	0, 2
	1	1, 0	1, 1	1, 2
	2	2, 0	2, 1	2, 2

In [17]:

```
1"""
2Higher dimensional numpy arrays give you more options,
3as you can slice one or more axes and also mix integers.
4Consider the 2D array above, arr2d. Slicing this array is a bit different
5"""
6print("Sliced array (matrix):\n{}".format(arr2d[:2]))
7"""
8As you can see, it has sliced along axis 0, the first axis.
9A slice, therefore, selects a range of elements along an axis.
10You can pass multiple slices just like you can pass multiple indexes
11"""
12# Extract the first two elements along axis 0 (i.e., the first two rows)
13# and every element except the first along axis 1 (i.e., the second and third columns)
14print("Sliced array (matrix):\n{}".format(arr2d[:2, 1:]))
15# When slicing like this, you always obtain array views of the same number of dimensions.
16# By mixing integer indexes and slices, you get lower dimensional slices
17# Access the whole second element along axis 0 and the first two along axis 1.
18print("Sliced array (matrix):\n{}".format(arr2d[1, :2]))
```

```
Sliced array (matrix):
[[1 2 3]
 [4 5 6]]
Sliced array (matrix):
[[2 3]
 [5 6]]
Sliced array (matrix):
[4 5]
```



In [18]:

```
1 """
2 Note that a colon by itself means to take the entire axis,
3 so you can slice only higher dimensional axes by doing as follows
4 """
5 # Extract the first column
6 print("Sliced array (matrix):\n{}".format(arr2d[:, :1]))
7 # Of course, assigning to a slice expression assigns to the whole selection
8 arr2d[:2, 1:] = 0
9 print("New array (matrix):\n{}".format(arr2d))
```

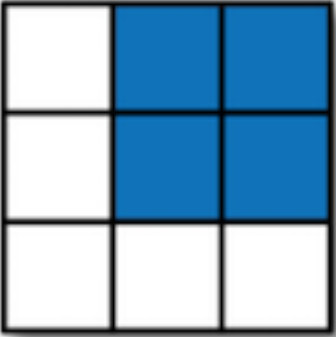
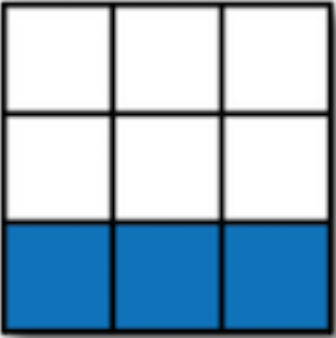
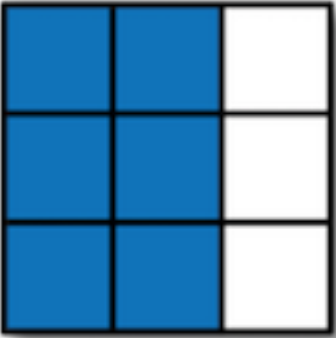
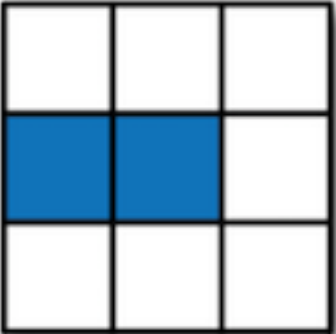
Sliced array (matrix):

```
[[1]
 [4]
 [7]]
```

New array (matrix):

```
[[1 0 0]
 [4 0 0]
 [7 8 9]]
```

# Slicing on a 2-d Array

	Expression	Shape
	<code>arr[:2, 1:]</code>	<code>(2, 2)</code>
	<code>arr[2]</code> <code>arr[2, :]</code> <code>arr[2:, :]</code>	<code>(3,)</code> <code>(3,)</code> <code>(1, 3)</code>
	<code>arr[:, :2]</code>	<code>(3, 2)</code>
	<code>arr[1, :2]</code> <code>arr[1:2, :2]</code>	<code>(2,)</code> <code>(1, 2)</code>

# Boolean Indexing

In [19]:

```
1"""
2Let's consider an array containing some data and an array of names with duplicates.
3We generate some random normally distributed data with the 'randn' function in numpy.random
4"""
5# Random, normally distributed 7x4 data matrix
6data = np.random.randn(7, 4)
7print("The original input data is:\n{}".format(data))
8# numpy array containing "names"
9names = np.array(['Bob', 'Joe', 'Will', 'Bob', 'Will', 'Joe', 'Joe'])
10"""
11Suppose each name corresponds to a row in the data array,
12and we want to select all the rows with corresponding name 'Bob'.
13Like arithmetic operations, comparisons (such as ==) with arrays are also vectorized.
14Thus, comparing names with the string 'Bob' yields a boolean array
15"""
16names == 'Bob'
```

The original input data is:

```
[[-0.7719735  -0.54500994  1.08807583  1.49070231]
 [-0.98309067 -0.70861066 -1.93274927 -1.79011088]
 [ 1.63070203 -0.53315322  0.58725735  0.17322339]
 [-0.9051374  -0.6947978  -1.96765755  0.53994032]
 [-0.24636087  0.33501854  0.71853993  0.89458756]
 [-1.23836331 -1.69812044  0.50381849 -1.61415737]
 [ 0.7691585   0.8033165   2.0940385  -0.49227147]]
```

Out[19]: array([ True, False, False, True, False, False, False], dtype=bool)

In [20]:

```
1"""
2The boolean array above can be passed when indexing the array
3"""
4data[names == 'Bob']
```

Out[20]: array([[ -0.7719735 , -0.54500994, 1.08807583, 1.49070231],  
 [ -0.9051374 , -0.6947978 , -1.96765755, 0.53994032]])

In [21]:

```
1"""
2The boolean array must be of the same length as the axis it is indexing.
3You can even mix and match boolean arrays with slices or integers (or sequences of integers)
4"""
5# Extract all the rows indexed by the boolean array yet limited to 3rd and 4th columns
6print("Boolean indexing, 3rd and 4th columns only:\n{}".format(data[names == 'Bob', 2:]))
7# Extract all the rows indexed by the boolean array yet limited to 2nd column
8print("Boolean indexing, 2nd column only:\n{}".format(data[names == 'Bob', 1]))
```

Boolean indexing, 3rd and 4th columns only:

```
[[ 1.08807583  1.49070231]
 [-1.96765755  0.53994032]]
```

Boolean indexing, 2nd column only:

```
[-0.54500994 -0.6947978 ]
```

In [22]:

```
1"""  
2To select everything but 'Bob', you can either use '!=' or negate the condition using '~'  
3"""  
4data[~(names == 'Bob')]
```

Out[22]: array([[ -0.98309067, -0.70861066, -1.93274927, -1.79011088],  
 [ 1.63070203, -0.53315322, 0.58725735, 0.17322339],  
 [-0.24636087, 0.33501854, 0.71853993, 0.89458756],  
 [-1.23836331, -1.69812044, 0.50381849, -1.61415737],  
 [ 0.7691585 , 0.8033165 , 2.0940385 , -0.49227147]])

In [23]:

```
1 """
2 To select more than one names to combine multiple boolean conditions,
3 use boolean arithmetic operators like '&' (and) and '|' (or)
4 NOTE: The Python keywords 'and' and 'or' DO NOT work with boolean arrays!!!
5 Selecting data from an array by boolean indexing always creates a copy of the data,
6 even if the returned array is unchanged.
7 """
8 mask = (names == 'Bob') | (names == 'Will')
9 print("Masked data:\n{}".format(data[mask]))
```

Masked data:

```
[[-0.7719735  -0.54500994  1.08807583  1.49070231]
 [ 1.63070203 -0.53315322  0.58725735  0.17322339]
 [-0.9051374  -0.6947978  -1.96765755  0.53994032]
 [-0.24636087  0.33501854  0.71853993  0.89458756]]
```



In [24]:

```
1 """
2 Setting values with boolean arrays works in a common-sense way.
3 To set all of the negative values in 'data' to '0' we need only to do the following.
4 """
5 data[data < 0] = 0
6 data
```

Out[24]: array([[ 0. , 0. , 1.08807583, 1.49070231],  
[ 0. , 0. , 0. , 0. ],  
[ 1.63070203, 0. , 0.58725735, 0.17322339],  
[ 0. , 0. , 0. , 0.53994032],  
[ 0. , 0.33501854, 0.71853993, 0.89458756],  
[ 0. , 0. , 0.50381849, 0. ],  
[ 0.7691585 , 0.8033165 , 2.0940385 , 0. ]])

In [25]:

```
1 """
2 Setting whole rows or columns using a 1D boolean array is also easy.
3 """
4 data[names != 'Joe'] = 5
5 data
```

Out[25]: array([[ 5. , 5. , 5. , 5. ],  
[ 0. , 0. , 0. , 0. ],  
[ 5. , 5. , 5. , 5. ],  
[ 5. , 5. , 5. , 5. ],  
[ 5. , 5. , 5. , 5. ],  
[ 0. , 0. , 0.50381849, 0. ],  
[ 0.7691585 , 0.8033165 , 2.0940385 , 0. ]])

# Transposing Arrays and Swapping Axes

In [26]:

```
1 """
2 Transposing is a special form of reshaping which returns a view on the underlying data
3 without copying anything.
4 Arrays have the transpose method and also the special 'T' attribute.
5 """
6 # Let's define a 1-d numpy array
7 arr = np.arange(15)
8 print("The original numpy array is: {}".format(arr))
9 reshaped_arr = arr.reshape((3, 5))
10 print("The reshaped numpy array is:\n{}".format(reshaped_arr))
11 transposed_arr = reshaped_arr.T
12 print("The transposed numpy array is:\n{}".format(transposed_arr))
```

The original numpy array is: [ 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14]

The reshaped numpy array is:

```
[[ 0  1  2  3  4]
 [ 5  6  7  8  9]
 [10 11 12 13 14]]
```

The transposed numpy array is:

```
[[ 0  5 10]
 [ 1  6 11]
 [ 2  7 12]
 [ 3  8 13]
 [ 4  9 14]]
```

In [27]:

```
1"""
2When doing matrix computations, you will do this very often,
3like for example computing the dot matrix product XTX using np.dot
4"""
5# Create a 4x3 matrix
6matrix = np.random.randn(4, 3)
7print("The original matrix is:\n{}".format(matrix))
8dot_product = np.dot(matrix.T, matrix)
9print("The result of the dot product is:\n{}".format(dot_product))
```

The original matrix is:

```
[[ 0.53942833 -0.60514056  0.17707711]
 [ 1.35496875 -0.61242005  1.15398291]
 [ 1.78242713  0.02099927  0.3451653 ]
 [-0.5397399   0.55667071  2.01884349]]
```

The result of the dot product is:

```
[[ 5.59528889 -1.41926772  1.18471281]
 [-1.41926772  1.05157666  0.31720044]
 [ 1.18471281  0.31720044  5.55790097]]
```

# Universal Functions: Fast Element-wise Array Functions

- A universal function, or **ufunc**, is a function that performs elementwise operations on data in **ndarrays**.
- You can think of them as fast vectorized wrappers for simple functions that take one or more scalar values and produce one or more scalar results.
- In case of binary universal functions, the shape of the input arrays **must be the same**.

In [28]:

```
1"""
2Many ufuncs are simple elementwise unary transformations, like 'sqrt' or 'exp'.
3"""
4arr = np.arange(10)
5print("The original array is: {}".format(arr))
6sqrt_arr = np.sqrt(arr)
7print("The squared-root array is: {}".format(sqrt_arr))
8exp_arr = np.exp(arr)
9print("The exp array is: {}".format(exp_arr))
10"""
11Other functions, such as 'add' or 'maximum', take 2 arrays (thus, binary ufuncs)
12and return a single array as the result
13"""
14# Define two random arrays
15x = np.random.randn(5)
16y = np.random.randn(5)
17print("x = {}".format(x))
18print("y = {}".format(y))
19print("Element-wise maximum between x's and y's elements: {}".format(np.maximum(x, y)))
```

The original array is: [0 1 2 3 4 5 6 7 8 9]

The squared-root array is: [ 0. 1. 1.41421356 1.73205081 2.  
2.23606798

2.44948974 2.64575131 2.82842712 3. ]

The exp array is: [ 1.00000000e+00 2.71828183e+00 7.38905610e+00 2.00855369e+01  
5.45981500e+01 1.48413159e+02 4.03428793e+02 1.09663316e+03  
2.98095799e+03 8.10308393e+03]

x = [-1.60788496 1.3027705 1.29542805 -0.5510918 1.55611235]

y = [-0.32853292 -0.36326418 0.19170485 0.15990394 0.76457901]

Element-wise maximum between x's and y's elements: [-0.32853292 1.3027705 1.29542805  
0.15990394 1.55611235]

# Universal Unary Functions (1 of 2)

Function	Description
<code>abs</code> , <code>fabs</code>	Compute the absolute value element-wise for integer, floating point, or complex values. Use <code>fabs</code> as a faster alternative for non-complex-valued data
<code>sqrt</code>	Compute the square root of each element. Equivalent to <code>arr ** 0.5</code>
<code>square</code>	Compute the square of each element. Equivalent to <code>arr ** 2</code>
<code>exp</code>	Compute the exponent $e^x$ of each element
<code>log</code> , <code>log10</code> , <code>log2</code> , <code>log1p</code>	Natural logarithm (base $e$ ), log base 10, log base 2, and $\log(1 + x)$ , respectively
<code>sign</code>	Compute the sign of each element: 1 (positive), 0 (zero), or -1 (negative)
<code>ceil</code>	Compute the ceiling of each element, i.e. the smallest integer greater than or equal to each element
<code>floor</code>	Compute the floor of each element, i.e. the largest integer less than or equal to each element
<code>rint</code>	Round elements to the nearest integer, preserving the dtype



# Universal Unary Functions (2 of 2)

<code>modf</code>	Return fractional and integral parts of array as separate array
<code>isnan</code>	Return boolean array indicating whether each value is NaN (Not a Number)
<code>isfinite, isinf</code>	Return boolean array indicating whether each element is finite (non-inf, non-NaN) or infinite, respectively
<code>cos, cosh, sin, sinh, tan, tanh</code>	Regular and hyperbolic trigonometric functions
<code>arccos, arccosh, arcsin, arcsinh, arctan, arctanh</code>	Inverse trigonometric functions
<code>logical_not</code>	Compute truth value of <code>not</code> x element-wise. Equivalent to <code>-arr</code> .

# Universal Binary Functions

Function	Description
<code>add</code>	Add corresponding elements in arrays
<code>subtract</code>	Subtract elements in second array from first array
<code>multiply</code>	Multiply array elements
<code>divide, floor_divide</code>	Divide or floor divide (truncating the remainder)
<code>power</code>	Raise elements in first array to powers indicated in second array
<code>maximum, fmax</code>	Element-wise maximum. <code>fmax</code> ignores NaN
<code>minimum, fmin</code>	Element-wise minimum. <code>fmin</code> ignores NaN
<code>mod</code>	Element-wise modulus (remainder of division)
<code>copysign</code>	Copy sign of values in second argument to values in first argument
<code>greater, greater_equal, less, less_equal, equal, not_equal</code>	Perform element-wise comparison, yielding boolean array. Equivalent to infix operators <code>&gt;</code> , <code>&gt;=</code> , <code>&lt;</code> , <code>&lt;=</code> , <code>==</code> , <code>!=</code>
<code>logical_and, logical_or, logical_xor</code>	Compute element-wise truth value of logical operation. Equivalent to infix operators <code>&amp;</code> , <code> </code> , <code>^</code>

# Mathematical and Statistical Methods

- A set of mathematical functions which compute statistics about an entire array or about the data along an axis are accessible as array methods.
- Aggregations (often called reductions) like `sum`, `mean`, and `std` (standard deviation) can be invoked:
  - by calling the array instance method;
  - using the top level `numpy` function.

In [29]:

```
# Consider the following normally-distributed random 5x4 matrix data
matrix = np.random.randn(5, 4)
print("The original matrix is:\n{}".format(matrix))
print("The mean of the matrix is: {}".format(matrix.mean()))
print("The mean of the matrix is: {}".format(np.mean(matrix)))
"""
Functions like 'mean' and 'sum' take an optional axis argument,
which computes the statistic over the given axis,
resulting in an array with one fewer dimension
"""
print("The mean of the matrix along the columns is: {}".format(matrix.mean(axis=1)))
print("The sum of the matrix along the rows is: {}".format(matrix.sum(axis=0)))
```

The original matrix is:

```
[[ 0.80315271  0.43595353  0.25475635 -1.41916843]
 [ 0.08995756  0.24721462  2.36291207  0.36666885]
 [ 0.24413369 -0.2324963   0.709209    0.05656641]
 [-1.49674819  0.25793226 -1.04658819  0.23984137]
 [-1.52227821  0.13571985  0.81259401  0.86914029]]
```

The mean of the matrix is: 0.1084236629677177

The mean of the matrix is: 0.1084236629677177

The mean of the matrix along the columns is: [ 0.01867354 0.76668828 0.1943532 -0.51139069 0.07379399]

The sum of the matrix along the rows is: [-1.88178244 0.84432396 3.09288324 0.1130485 ]

# Table of `numpy` Statistical Methods

Method	Description
<code>sum</code>	Sum of all the elements in the array or along an axis. Zero-length arrays have sum 0.
<code>mean</code>	Arithmetic mean. Zero-length arrays have NaN mean.
<code>std</code> , <code>var</code>	Standard deviation and variance, respectively, with optional degrees of freedom adjustment (default denominator n).
<code>min</code> , <code>max</code>	Minimum and maximum.
<code>argmin</code> , <code>argmax</code>	Indices of minimum and maximum elements, respectively.
<code>cumsum</code>	Cumulative sum of elements starting from 0
<code>cumprod</code>	Cumulative product of elements starting from 1

# Table of `numpy` Set Methods

Method	Description
<code>unique(x)</code>	Compute the sorted, unique elements in <code>x</code>
<code>intersect1d(x, y)</code>	Compute the sorted, common elements in <code>x</code> and <code>y</code>
<code>union1d(x, y)</code>	Compute the sorted union of elements
<code>in1d(x, y)</code>	Compute a boolean array indicating whether each element of <code>x</code> is contained in <code>y</code>
<code>setdiff1d(x, y)</code>	Set difference, elements in <code>x</code> that are not in <code>y</code>
<code>setxor1d(x, y)</code>	Set symmetric differences; elements that are in either of the arrays, but not both

# I/O with `numpy` Arrays

- `numpy` is able to **save** and **load** data to and from disk either in *text* or *binary* format.
- We only discuss built-in binary format, since we will use `pandas` for loading text or tabular data.
- `np.save` and `np.load` are the two workhorse functions for efficiently saving and loading array data on disk.
- Arrays are saved by default in an *uncompressed raw binary* format with file extension `.npy`



In [30]:

```
1# Consider the following numpy array
arr = np.random.randn(5)
print("The original array is: {}".format(arr))
4# Persist the above array out to disk to the specified path on disk
np.save("./data/np_array", arr) # if no '.npy' extension is specified it will be appended
6# Load the array back from the specified path on disk
arr_loaded = np.load("./data/np_array.npy")
print("The array loaded from disk is: {}".format(arr_loaded))
9
1# NOTE: If you need to save multiple arrays in a zip archive
1# use 'np.savez' and pass the arrays as keyword arguments:
1#     np.savez("path/to/arr_archive.npz", a=arr_a, b=arr_b)
1# When this is loaded back with:
1#     arr_archive = np.load("./data/np_array.npy")
1# You get back a dict-like object which loads the individual arrays lazily:
1#     arr_archive['b'] # refers to the second array in the archive
```

The original array is: [-0.00948631 -0.18875794 -0.49668335 -0.08650266 1.06290762]

The array loaded from disk is: [-0.00948631 -0.18875794 -0.49668335 -0.08650266 1.0629076

2]



# Table of `numpy` Linear Algebra Functions

Function	Description
<code>diag</code>	Return the diagonal (or off-diagonal) elements of a square matrix as a 1D array, or convert a 1D array into a square matrix with zeros on the off-diagonal
<code>dot</code>	Matrix multiplication
<code>trace</code>	Compute the sum of the diagonal elements
<code>det</code>	Compute the matrix determinant
<code>eig</code>	Compute the eigenvalues and eigenvectors of a square matrix
<code>inv</code>	Compute the inverse of a square matrix
<code>pinv</code>	Compute the Moore-Penrose pseudo-inverse inverse of a matrix
<code>qr</code>	Compute the QR decomposition
<code>svd</code>	Compute the singular value decomposition (SVD)
<code>solve</code>	Solve the linear system $Ax = b$ for $x$ , where $A$ is a square matrix
<code>lstsq</code>	Compute the least-squares solution to $Ax = b$

# Table of `numpy.random` Functions

Function	Description
<code>seed</code>	Seed the random number generator
<code>permutation</code>	Return a random permutation of a sequence, or return a permuted range
<code>shuffle</code>	Randomly permute a sequence in place
<code>rand</code>	Draw samples from a uniform distribution
<code>randint</code>	Draw random integers from a given low-to-high range
<code>randn</code>	Draw samples from a normal distribution with mean 0 and standard deviation 1 (MATLAB-like interface)
<code>binomial</code>	Draw samples from a binomial distribution
<code>normal</code>	Draw samples from a normal (Gaussian) distribution
<code>beta</code>	Draw samples from a beta distribution
<code>chisquare</code>	Draw samples from a chi-square distribution
<code>gamma</code>	Draw samples from a gamma distribution
<code>uniform</code>	Draw samples from a uniform [0, 1) distribution