

# 01b - Introduction - Scientific Python

January 15, 2017

## 1 Python for scientific computing

Python has extensive packages to help with data analysis:

- numpy: matrices, linear algebra, Fourier transform, pseudorandom number generators
- scipy: advanced linear algebra and maths, signal processing, statistics
- pandas: DataFrames, data wrangling and analysis
- matplotlib: visualizations such as line charts, histograms, scatter plots.

```
In [96]: from preamble import * # This is to simplify the code.
         %matplotlib inline
         HTML('<style>.CodeMirror{min-width:100% !important;}</style>') # For slides
```

```
Out[96]: <IPython.core.display.HTML object>
```

### 1.1 NumPy

NumPy is the fundamental package required for high performance scientific computing in Python. It provides:

- ndarray: fast and space-efficient n-dimensional numeric array with vectorized arithmetic operations
- Functions for fast operations on arrays without having to write loops
- Linear algebra, random number generation, Fourier transform
- Integrating code written in C, C++, and Fortran (for faster operations)

pandas provides a richer, simpler interface to many operations. We'll focus on using ndarrays here because they are heavily used in scikit-learn.

#### 1.1.1 ndarrays

There are several ways to create numpy arrays.

```
In [97]: # Convert normal Python array to 1-dimensional numpy array
         np.array((1, 2, 53))
```

```
Out[97]: array([ 1,  2, 53])
```

```
In [98]: # Convert sequences of sequences of sequences ... to n-dim array
np.array([(1.5, 2, 3), (4, 5, 6)])
```

```
Out[98]: array([[ 1.5,  2. ,  3. ],
                [ 4. ,  5. ,  6. ]])
```

```
In [99]: # Define element type at creation time
np.array([[1, 2], [3, 4]], dtype=complex)
```

```
Out[99]: array([[ 1.+0.j,  2.+0.j],
                [ 3.+0.j,  4.+0.j]])
```

Useful properties of ndarrays:

```
In [100]: my_array = np.array([[1, 0, 3], [0, 1, 2]])
my_array.ndim      # number of dimensions (axes), also called the rank
my_array.shape     # a matrix with n rows and m columns has shape (n,m)
my_array.size      # the total number of elements of the array
my_array.dtype     # type of the elements in the array
my_array.itemsize  # the size in bytes of each element of the array
```

```
Out[100]: 2
```

```
Out[100]: (2, 3)
```

```
Out[100]: 6
```

```
Out[100]: dtype('int64')
```

```
Out[100]: 8
```

Quick array creation.

It is cheaper to create an array with placeholders than extending it later.

```
In [101]: np.ones(3) # Default type is float64
np.zeros([2, 2])
np.empty([2, 2]) # Fills the array with whatever sits in memory
np.random.random((2,3))
np.random.randint(5, size=(2, 4))
```

```
Out[101]: array([ 1.,  1.,  1.])
```

```
Out[101]: array([[ 0.,  0.],
                 [ 0.,  0.]])
```

```
Out[101]: array([[ 0.,  0.],
                 [ 0.,  0.]])
```

```
Out[101]: array([[ 0.441,  0.933,  0.906],
                 [ 0.741,  0.681,  0.462]])
```

```
Out[101]: array([[3, 3, 3, 2],
                [2, 1, 3, 4]])
```

Create sequences of numbers

```
In [102]: np.linspace(0, 1, num=4) # Linearly distributed numbers between 0 and 1
          np.arange(0, 1, step=0.3) # Fixed step size
          np.arange(12).reshape(3,4) # Create and reshape
          np.eye(4) # Identity matrix
```

```
Out[102]: array([ 0.    ,  0.333,  0.667,  1.    ])
```

```
Out[102]: array([ 0. ,  0.3,  0.6,  0.9])
```

```
Out[102]: array([[ 0,  1,  2,  3],
                [ 4,  5,  6,  7],
                [ 8,  9, 10, 11]])
```

```
Out[102]: array([[ 1.,  0.,  0.,  0.],
                [ 0.,  1.,  0.,  0.],
                [ 0.,  0.,  1.,  0.],
                [ 0.,  0.,  0.,  1.]])
```

## 1.1.2 Basic Operations

Arithmetic operators on arrays apply elementwise. A new array is created and filled with the result. Some operations, such as += and \*=, act in place to modify an existing array rather than create a new one.

```
In [103]: a = np.array([20, 30, 40, 50])
          b = np.arange(4)
          a, b # Just printing
          a-b
          b**2
          a > 32
          a += 1
          a
```

```
Out[103]: (array([20, 30, 40, 50]), array([0, 1, 2, 3]))
```

```
Out[103]: array([20, 29, 38, 47])
```

```
Out[103]: array([0, 1, 4, 9])
```

```
Out[103]: array([False, False,  True,  True], dtype=bool)
```

```
Out[103]: array([21, 31, 41, 51])
```

The product operator \* operates elementwise. The matrix product can be performed using dot()

```
In [104]: A, B = np.array([[1,1], [0,1]]), np.array([[2,0], [3,4]]) # assign multiple variables
A
B
A * B
np.dot(A, B)
```

```
Out[104]: array([[1, 1],
                [0, 1]])
```

```
Out[104]: array([[2, 0],
                [3, 4]])
```

```
Out[104]: array([[2, 0],
                [0, 4]])
```

```
Out[104]: array([[5, 4],
                [3, 4]])
```

Upcasting: Operations with arrays of different types choose the more general/precise one.

```
In [105]: a = np.ones(3, dtype=np.int) # initialize to integers
b = np.linspace(0, np.pi, 3) # default type is float
a.dtype, b.dtype, (a + b).dtype
```

```
Out[105]: (dtype('int64'), dtype('float64'), dtype('float64'))
```

ndarrays have most unary operations (max,min,sum,...) built in

```
In [106]: a = np.random.random((2,3))
a
a.sum(), a.min(), a.max()
```

```
Out[106]: array([[ 0.721,  0.43 ,  0.016],
                [ 0.287,  0.609,  0.247]])
```

```
Out[106]: (2.3097201412037407, 0.016272921833643483, 0.72108163821790305)
```

By specifying the axis parameter you can apply an operation along a specified axis of an array

```
In [107]: b = np.arange(12).reshape(3,4)
b
b.sum(axis=0)
b.sum(axis=1)
```

```
Out[107]: array([[ 0,  1,  2,  3],
                [ 4,  5,  6,  7],
                [ 8,  9, 10, 11]])
```

```
Out[107]: array([12, 15, 18, 21])
```

```
Out[107]: array([ 6, 22, 38])
```

### 1.1.3 Universal Functions

NumPy provides familiar mathematical functions such as `sin`, `cos`, `exp`, `sqrt`, `floor`,... In NumPy, these are called "universal functions" (ufunc), and operate elementwise on an array, producing an array as output.

```
In [108]: np.sqrt(np.arange(0, 10))
```

```
Out[108]: array([ 0.    ,  1.    ,  1.414,  1.732,  2.    ,  2.236,  2.449,  2.646,
                2.828,  3.    ])
```

### 1.1.4 Shape Manipulation

Transpose, flatten, reshape,...

```
In [109]: a = np.floor(10*np.random.random((3,4)))
```

```
a
```

```
a.transpose()
```

```
b = a.ravel() # flatten array
```

```
b
```

```
b.reshape(3, -1) # reshape in 2 rows (and as many columns as needed)
```

```
Out[109]: array([[ 4.,  9.,  8.,  1.],
                 [ 7.,  2.,  2.,  6.],
                 [ 3.,  9.,  4.,  3.]])
```

```
Out[109]: array([[ 4.,  7.,  3.],
                 [ 9.,  2.,  9.],
                 [ 8.,  2.,  4.],
                 [ 1.,  6.,  3.]])
```

```
Out[109]: array([ 4.,  9.,  8.,  1.,  7.,  2.,  2.,  6.,  3.,  9.,  4.,  3.])
```

```
Out[109]: array([[ 4.,  9.,  8.,  1.],
                 [ 7.,  2.,  2.,  6.],
                 [ 3.,  9.,  4.,  3.]])
```

Arrays can be split and stacked together

```
In [110]: a = np.floor(10*np.random.random((2,6)))
```

```
a
```

```
b, c = np.hsplit(a, 2) # Idem: vsplit for vertical splits
```

```
b
```

```
c
```

```
np.hstack((b, c)) # Idem: vstack for vertical stacks
```

```
Out[110]: array([[ 9.,  5.,  5.,  1.,  8.,  1.],
                 [ 5.,  4.,  6.,  5.,  2.,  7.]])
```

```
Out[110]: array([[ 9.,  5.,  5.],
                 [ 5.,  4.,  6.]])
```

```
Out[110]: array([[ 1.,  8.,  1.],
                [ 5.,  2.,  7.]])
```

```
Out[110]: array([[ 9.,  5.,  5.,  1.,  8.,  1.],
                [ 5.,  4.,  6.,  5.,  2.,  7.]])
```

### 1.1.5 Indexing and Slicing

Arrays can be indexed and sliced using [start:stop:stepsize]. Defaults are [0:ndim:1]

```
In [111]: a = np.arange(10)**2
          a
```

```
Out[111]: array([ 0,  1,  4,  9, 16, 25, 36, 49, 64, 81])
```

```
In [112]: a[2]
```

```
Out[112]: 4
```

```
In [113]: a[3:10:2]
```

```
Out[113]: array([ 9, 25, 49, 81])
```

```
In [114]: a[::-1] # Defaults are used if indices not stated
```

```
Out[114]: array([81, 64, 49, 36, 25, 16,  9,  4,  1,  0])
```

```
In [115]: a[::2]
```

```
Out[115]: array([ 0,  4, 16, 36, 64])
```

For multi-dimensional arrays, axes are comma-separated: [x,y,z].

```
In [116]: b = np.arange(16).reshape(4,4)
          b
          b[2,3] # row 2, column 3
```

```
Out[116]: array([[ 0,  1,  2,  3],
                [ 4,  5,  6,  7],
                [ 8,  9, 10, 11],
                [12, 13, 14, 15]])
```

```
Out[116]: 11
```

```
In [117]: b[0:3,1] # Values 0 to 3 in column 1
          b[ : ,1] # The whole column 1
```

```
Out[117]: array([1, 5, 9])
```

```
Out[117]: array([ 1,  5,  9, 13])
```

```
In [118]: b[1:3, :] # Rows 1:3, all columns
```

```
Out[118]: array([[ 4,  5,  6,  7],
                [ 8,  9, 10, 11]])
```

```
In [119]: # Return the last row
```

```
b[-1]
```

```
Out[119]: array([12, 13, 14, 15])
```

Note: dots (...) represent as many colons (:) as needed \*  $x[1,2,\dots] = x[1,2,:::]$  \*  $x[\dots,3] = x[:,:::,3]$   
\*  $x[4,\dots,5,:] = x[4,::,5,:]$

Arrays can also be indexed by arrays of integers and booleans.

```
In [120]: a = np.arange(12)**2
          i = np.array([ 1,1,3,8,5 ])
          a
          a[i]
```

```
Out[120]: array([ 0,  1,  4,  9, 16, 25, 36, 49, 64, 81, 100, 121])
```

```
Out[120]: array([ 1,  1,  9, 64, 25])
```

A matrix of indices returns a matrix with the corresponding values.

```
In [121]: j = np.array([[ 3, 4], [9, 7]])
          a[j]
```

```
Out[121]: array([[ 9, 16],
                [81, 49]])
```

With boolean indices we explicitly choose which items in the array we want and which ones we don't.

```
In [122]: a = np.arange(12).reshape(3,4)
          a
          a[np.array([False,True,True]), :]
          b = a > 4
          b
          a[b]
```

```
Out[122]: array([[ 0,  1,  2,  3],
                [ 4,  5,  6,  7],
                [ 8,  9, 10, 11]])
```

```
Out[122]: array([[ 4,  5,  6,  7],
                [ 8,  9, 10, 11]])
```

```
Out[122]: array([[False, False, False, False],
                [False,  True,  True,  True],
                [ True,  True,  True,  True]], dtype=bool)
```

```
Out[122]: array([ 5,  6,  7,  8,  9, 10, 11])
```

### 1.1.6 Iterating

Iterating is done with respect to the first axis:

```
In [123]: for row in b:
           print(row)

[False False False False]
[False True True True]
[ True True True True]
```

Operations on each element can be done by flattening the array (or nested loops)

```
In [124]: for element in b.flat: # flat returns an iterator
           print(element)

False
False
False
False
False
True
True
True
True
True
True
True
True
```

### 1.1.7 Copies and Views (or: how to shoot yourself in a foot)

Assigning an array to another variable does NOT create a copy

```
In [125]: a = np.arange(12)
           b = a
           a
```

```
Out[125]: array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11])
```

```
In [126]: b[0] = -100
           b
```

```
Out[126]: array([-100,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10,
                11])
```

```
In [127]: a
```

```
Out[127]: array([-100,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10,
                11])
```



The `view()` method creates a NEW array object that looks at the same data.

```
In [128]: a = np.arange(12)
          a
          c = a.view()
          c.resize((2, 6))
          c
```

```
Out[128]: array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11])
```

```
Out[128]: array([[ 0,  1,  2,  3,  4,  5],
                 [ 6,  7,  8,  9, 10, 11]])
```

```
In [129]: a[0] = 123
          c # c is also changed now
```

```
Out[129]: array([[123,  1,  2,  3,  4,  5],
                 [ 6,  7,  8,  9, 10, 11]])
```

Slicing an array returns a view of it.

```
In [130]: c
          s = c[ : , 1:3]
          s[:] = 10
          s
          c
```

```
Out[130]: array([[123,  1,  2,  3,  4,  5],
                 [ 6,  7,  8,  9, 10, 11]])
```

```
Out[130]: array([[10, 10],
                 [10, 10]])
```

```
Out[130]: array([[123, 10, 10,  3,  4,  5],
                 [ 6, 10, 10,  9, 10, 11]])
```

The `copy()` method makes a deep copy of the array and its data.

```
In [131]: d = a.copy()
          d[0] = -42
          d
```

```
Out[131]: array([-42, 10, 10,  3,  4,  5,  6, 10, 10,  9, 10, 11])
```

```
In [132]: a
```

```
Out[132]: array([123, 10, 10,  3,  4,  5,  6, 10, 10,  9, 10, 11])
```

### 1.1.8 Numpy: further reading

- Numpy Tutorial: [http://wiki.scipy.org/Tentative\\_NumPy\\_Tutorial](http://wiki.scipy.org/Tentative_NumPy_Tutorial)
- "Python for Data Analysis" by Wes McKinney (O'Reilly)

## 1.2 SciPy

SciPy is a collection of packages for scientific computing, among others:

- `scipy.integrate`: numerical integration and differential equation solvers
- `scipy.linalg`: linear algebra routines and matrix decompositions
- `scipy.optimize`: function optimizers (minimizers) and root finding algorithms
- `scipy.signal`: signal processing tools
- `scipy.sparse`: sparse matrices and sparse linear system solvers
- `scipy.stats`: probability distributions, statistical tests, descriptive statistics

### 1.2.1 Sparse matrices

Sparse matrices are used in scikit-learn for (large) arrays that contain mostly zeros. You can convert a dense (numpy) matrix to a sparse matrix.

```
In [133]: from scipy import sparse
          eye = np.eye(4)
          eye
          sparse_matrix = sparse.csr_matrix(eye) # Compressed Sparse Row matrix
          sparse_matrix
          print("{}".format(sparse_matrix))
```

```
Out[133]: array([[ 1.,  0.,  0.,  0.],
                 [ 0.,  1.,  0.,  0.],
                 [ 0.,  0.,  1.,  0.],
                 [ 0.,  0.,  0.,  1.]])
```

```
Out[133]: <4x4 sparse matrix of type '<class 'numpy.float64'>'
          with 4 stored elements in Compressed Sparse Row format>
```

```
(0, 0)      1.0
(1, 1)      1.0
(2, 2)      1.0
(3, 3)      1.0
```

When the data is too large, you can create a sparse matrix by passing the values and coordinates (COO format).

```
In [134]: data = np.ones(4) # [1,1,1,1]
          row_indices = col_indices = np.arange(4) # [1,2,3,4]
          eye_coo = sparse.coo_matrix((data, (row_indices, col_indices)))
          print("{}".format(eye_coo))
```

```
(0, 0)      1.0
(1, 1)      1.0
(2, 2)      1.0
(3, 3)      1.0
```

## 1.2.2 Further reading

Check the [SciPy reference guide](#) for tutorials and examples of all SciPy capabilities.

## 1.3 pandas

pandas is a Python library for data wrangling and analysis. It provides:

- DataFrame: a table, similar to an R DataFrame that holds any structured data
  - Every column can have its own data type (strings, dates, floats,...)
- A great range of methods to apply to this table (sorting, querying, joining,...)
- Imports data from a wide range of data formats (CSV, Excel) and databases (e.g. SQL)

### 1.3.1 Series

A one-dimensional array of data (of any numpy type), with indexed values. It can be created by passing a Python list or dict, a numpy array, a csv file,...

```
In [135]: import pandas as pd
          pd.Series([1,3,np.nan]) # Default integers are integers
          pd.Series([1,3,5], index=['a','b','c'])
          pd.Series({'a' : 1, 'b': 2, 'c': 3 }) # when given a dict, the keys will be used for t
          pd.Series({'a' : 1, 'b': 2, 'c': 3 }, index = ['b', 'c', 'd']) # this will try to matc
```

```
Out [135]: 0    1.0
           1    3.0
           2    NaN
           dtype: float64
```

```
Out [135]: a    1
           b    3
           c    5
           dtype: int64
```

```
Out [135]: a    1
           b    2
           c    3
           dtype: int64
```

```
Out [135]: b    2.0
           c    3.0
           d    NaN
           dtype: float64
```

Functions like a numpy array, however with index labels as indices

```
In [136]: a = pd.Series({'a' : 1, 'b': 2, 'c': 3 })
          a
          a['b'] # Retrieves a value
          a[['a','b']] # and can also be sliced
```

```
Out[136]: a    1
          b    2
          c    3
          dtype: int64
```

```
Out[136]: 2
```

```
Out[136]: a    1
          b    2
          dtype: int64
```

numpy array operations on Series preserve the index value

```
In [137]: a
          a[a > 1]
          a * 2
          np.sqrt(a)
```

```
Out[137]: a    1
          b    2
          c    3
          dtype: int64
```

```
Out[137]: b    2
          c    3
          dtype: int64
```

```
Out[137]: a    2
          b    4
          c    6
          dtype: int64
```

```
Out[137]: a    1.00
          b    1.41
          c    1.73
          dtype: float64
```

Operations over multiple Series will align the indices

```
In [138]: a = pd.Series({'John' : 1000, 'Mary': 2000, 'Andre': 3000 })
          b = pd.Series({'John' : 100, 'Andre': 200, 'Cecilia': 300 })
          a + b
```

```
Out[138]: Andre      3200.0
          Cecilia      NaN
          John        1100.0
          Mary         NaN
          dtype: float64
```

### 1.3.2 DataFrame

A DataFrame is a tabular data structure with both a row and a column index. It can be created by passing a dict of arrays, a csv file,...

```
In [139]: data = {'state': ['Ohio', 'Ohio', 'Nevada', 'Nevada'], 'year': [2000, 2001, 2001, 2002],
                'pop': [1.5, 1.7, 2.4, 2.9]}
pd.DataFrame(data)
pd.DataFrame(data, columns=['year', 'state', 'pop', 'color']) # Will match indices
```

```
Out[139]:
```

	pop	state	year
0	1.5	Ohio	2000
1	1.7	Ohio	2001
2	2.4	Nevada	2001
3	2.9	Nevada	2002

```
Out[139]:
```

	year	state	pop	color
0	2000	Ohio	1.5	NaN
1	2001	Ohio	1.7	NaN
2	2001	Nevada	2.4	NaN
3	2002	Nevada	2.9	NaN

It can be composed with a numpy array and row and column indices, and decomposed

```
In [140]: dates = pd.date_range('20130101', periods=4)
df = pd.DataFrame(np.random.randn(4,4), index=dates, columns=list('ABCD'))
df
```

```
Out[140]:
```

	A	B	C	D
2013-01-01	-0.14	-0.07	-1.21	0.03
2013-01-02	-0.94	0.57	0.86	1.24
2013-01-03	-0.04	0.46	-1.09	-0.31
2013-01-04	-1.72	0.56	-0.31	0.74

```
In [141]: df.index
df.columns
df.values
```

```
Out[141]: DatetimeIndex(['2013-01-01', '2013-01-02', '2013-01-03', '2013-01-04'], dtype='datetime64[ns]', freq='D')
```

```
Out[141]: Index(['A', 'B', 'C', 'D'], dtype='object')
```

```
Out[141]: array([[ -0.136, -0.072, -1.205,  0.033],
                 [-0.936,  0.571,  0.864,  1.236],
                 [-0.038,  0.459, -1.086, -0.307],
                 [-1.725,  0.561, -0.314,  0.742]])
```

DataFrames can easily read/write data from/to files

- `read_csv(source)`: load CSV data from file or url

- `read_table(source, sep=',')`: load delimited data with separator
- `df.to_csv(target)`: writes the DataFrame to a file

```
In [142]: dfs = pd.read_csv('data.csv')
dfs
dfs.set_value(0, 'a', 10)
dfs.to_csv('data.csv', index=False) # Don't export the row index
```

```
Out[142]:
```

	a	b	c	d	message
0	10	2	3	4	hello
1	5	6	7	8	world
2	9	10	11	12	foo

```
Out[142]:
```

	a	b	c	d	message
0	10	2	3	4	hello
1	5	6	7	8	world
2	9	10	11	12	foo

### 1.3.3 Simple operations

```
In [143]: df.head() # First 5 rows
df.tail() # Last 5 rows
```

```
Out[143]:
```

	A	B	C	D
2013-01-01	-0.14	-0.07	-1.21	0.03
2013-01-02	-0.94	0.57	0.86	1.24
2013-01-03	-0.04	0.46	-1.09	-0.31
2013-01-04	-1.72	0.56	-0.31	0.74

```
Out[143]:
```

	A	B	C	D
2013-01-01	-0.14	-0.07	-1.21	0.03
2013-01-02	-0.94	0.57	0.86	1.24
2013-01-03	-0.04	0.46	-1.09	-0.31
2013-01-04	-1.72	0.56	-0.31	0.74

```
In [144]: # Quick stats
df.describe()
```

```
Out[144]:
```

	A	B	C	D
count	4.00	4.00	4.00	4.00
mean	-0.71	0.38	-0.44	0.43
std	0.79	0.31	0.95	0.69
min	-1.72	-0.07	-1.21	-0.31
25%	-1.13	0.33	-1.12	-0.05
50%	-0.54	0.51	-0.70	0.39
75%	-0.11	0.56	-0.02	0.87
max	-0.04	0.57	0.86	1.24

```
In [145]: # Transpose
df.T
```

```
Out[145]:
```

	2013-01-01	2013-01-02	2013-01-03	2013-01-04
A	-0.14	-0.94	-0.04	-1.72
B	-0.07	0.57	0.46	0.56
C	-1.21	0.86	-1.09	-0.31
D	0.03	1.24	-0.31	0.74

```
In [146]: df.sort_index(axis=1, ascending=False) # Sort by index labels
df.sort(columns='B') # Sort by values
```

```
Out[146]:
```

	D	C	B	A
2013-01-01	0.03	-1.21	-0.07	-0.14
2013-01-02	1.24	0.86	0.57	-0.94
2013-01-03	-0.31	-1.09	0.46	-0.04
2013-01-04	0.74	-0.31	0.56	-1.72

```
Out[146]:
```

	A	B	C	D
2013-01-01	-0.14	-0.07	-1.21	0.03
2013-01-03	-0.04	0.46	-1.09	-0.31
2013-01-04	-1.72	0.56	-0.31	0.74
2013-01-02	-0.94	0.57	0.86	1.24

### 1.3.4 Selecting and slicing

```
In [147]: df['A'] # Get single column by label
df.A # Shorthand
```

```
Out[147]:
```

2013-01-01	-0.14
2013-01-02	-0.94
2013-01-03	-0.04
2013-01-04	-1.72

Freq: D, Name: A, dtype: float64

```
Out[147]:
```

2013-01-01	-0.14
2013-01-02	-0.94
2013-01-03	-0.04
2013-01-04	-1.72

Freq: D, Name: A, dtype: float64

```
In [148]: df[0:2] # Get rows by index number
df.iloc[0:2,0:2] # Get rows and columns by index number
df['20130102':'20130103'] # or row label
df.loc['20130102':'20130103', ['A','B']] # or row and column label
df.ix[0:2, ['A','B']] # allows mixing integers and labels
```

```
Out[148]:
```

	A	B	C	D
2013-01-01	-0.14	-0.07	-1.21	0.03
2013-01-02	-0.94	0.57	0.86	1.24

```
Out[148]:
```

	A	B
2013-01-01	-0.14	-0.07
2013-01-02	-0.94	0.57

```
Out [148]:
```

	A	B	C	D
2013-01-02	-0.94	0.57	0.86	1.24
2013-01-03	-0.04	0.46	-1.09	-0.31

```
Out [148]:
```

	A	B
2013-01-02	-0.94	0.57
2013-01-03	-0.04	0.46

```
Out [148]:
```

	A	B
2013-01-01	-0.14	-0.07
2013-01-02	-0.94	0.57

query() retrieves data matching a boolean expression

```
In [149]: df
df.query('A > 0.4') # Identical to df[df.A > 0.4]
df.query('A > B')  # Identical to df[df.A > df.B]
```

```
Out [149]:
```

	A	B	C	D
2013-01-01	-0.14	-0.07	-1.21	0.03
2013-01-02	-0.94	0.57	0.86	1.24
2013-01-03	-0.04	0.46	-1.09	-0.31
2013-01-04	-1.72	0.56	-0.31	0.74

```
Out [149]: Empty DataFrame
Columns: [A, B, C, D]
Index: []
```

```
Out [149]: Empty DataFrame
Columns: [A, B, C, D]
Index: []
```

Note: similar to NumPy, indexing and slicing returns a *view* on the data. Use copy() to make a deep copy.

### 1.3.5 Operations

DataFrames offer a [wide range of operations](#): max, mean, min, sum, std,...

```
In [150]: df.mean()          # Mean of all values per column
df.mean(axis=1) # Other axis: means per row
```

```
Out [150]: A    -0.71
           B     0.38
           C   -0.44
           D     0.43
           dtype: float64
```



```
Out [150]: 2013-01-01    -0.35
           2013-01-02     0.43
           2013-01-03   -0.24
           2013-01-04   -0.18
           Freq: D, dtype: float64
```

All of numpy's universal functions also work with dataframes

```
In [151]: np.abs(df)
```

```
Out [151]:
```

	A	B	C	D
2013-01-01	0.14	0.07	1.21	0.03
2013-01-02	0.94	0.57	0.86	1.24
2013-01-03	0.04	0.46	1.09	0.31
2013-01-04	1.72	0.56	0.31	0.74

Other (custom) functions can be applied with apply(func)

```
In [152]: df
           df.apply(np.max)
           df.apply(lambda x: x.max() - x.min())
```

```
Out [152]:
```

	A	B	C	D
2013-01-01	-0.14	-0.07	-1.21	0.03
2013-01-02	-0.94	0.57	0.86	1.24
2013-01-03	-0.04	0.46	-1.09	-0.31
2013-01-04	-1.72	0.56	-0.31	0.74

```
Out [152]: A    -0.04
           B     0.57
           C     0.86
           D     1.24
           dtype: float64
```

```
Out [152]: A     1.69
           B     0.64
           C     2.07
           D     1.54
           dtype: float64
```

Data can be aggregated with groupby()

```
In [153]: df = pd.DataFrame({'A' : ['foo', 'bar', 'foo', 'bar'], 'B' : ['one', 'one', 'two', 'th'],
                             'C' : np.random.randn(4), 'D' : np.random.randn(4)})
           df
           df.groupby('A').sum()
           df.groupby(['A', 'B']).sum()
```

```
Out[153]:
```

	A	B	C	D
0	foo	one	-1.38	-1.79
1	bar	one	-1.03	1.53
2	foo	two	-0.18	0.41
3	bar	three	1.04	0.19

```
Out[153]:
```

		C	D
A			
bar	6.68e-03	1.73	
foo	-1.56e+00	-1.39	

```
Out[153]:
```

	A	B	C	D
bar	one		-1.03	1.53
		three	1.04	0.19
foo	one		-1.38	-1.79
		two	-0.18	0.41

### 1.3.6 Data wrangling (some examples)

Merge: combine two dataframes based on common keys

```
In [154]: df1 = pd.DataFrame({'key': ['b', 'b', 'a'], 'data1': range(3)})
df2 = pd.DataFrame({'key': ['a', 'b'], 'data2': range(2)})
df1
df2
pd.merge(df1, df2)
```

```
Out[154]:
```

	data1	key
0	0	b
1	1	b
2	2	a

```
Out[154]:
```

	data2	key
0	0	a
1	1	b

```
Out[154]:
```

	data1	key	data2
0	0	b	1
1	1	b	1
2	2	a	0

Append: append one dataframe to another

```
In [155]: df = pd.DataFrame(np.random.randn(2, 4))
df
s = pd.DataFrame(np.random.randn(1,4))
s
df.append(s, ignore_index=True)
```

```
Out [155]:      0      1      2      3
0  1.71 -0.47  0.22  1.51
1 -1.57  0.27 -0.89  1.29
```

```
Out [155]:      0      1      2      3
0 -1.42  0.39 -1.57 -0.54
```

```
Out [155]:      0      1      2      3
0  1.71 -0.47  0.22  1.51
1 -1.57  0.27 -0.89  1.29
2 -1.42  0.39 -1.57 -0.54
```

### Remove duplicates

```
In [156]: df = pd.DataFrame({'k1': ['one'] * 3, 'k2': [1, 1, 2]})
df
df.drop_duplicates()
```

```
Out [156]:   k1  k2
0  one  1
1  one  1
2  one  2
```

```
Out [156]:   k1  k2
0  one  1
2  one  2
```

### Replace values

```
In [157]: df = pd.DataFrame({'k1': [1, -1], 'k2': [-1, 2]}) # Say that -1 is a sentinel for miss
df
df.replace(-1, np.nan)
```

```
Out [157]:   k1  k2
0    1 -1
1   -1  2
```

```
Out [157]:   k1  k2
0  1.0 NaN
1  NaN  2.0
```

### Discretization and binning

```
In [158]: ages = [20, 22, 25, 27, 21, 23, 37, 31, 61, 45, 41, 32]
bins = [18, 25, 35, 60, 100]
cats = pd.cut(ages, bins)
cats.labels
pd.value_counts(cats)
```

```
Out [158]: array([0, 0, 0, 1, 0, 0, 2, 1, 3, 2, 2, 1], dtype=int8)
```

```
Out[158]: (18, 25]    5
          (35, 60]    3
          (25, 35]    3
          (60, 100]   1
          dtype: int64
```

### 1.3.7 Further reading

- Pandas docs: <http://pandas.pydata.org/pandas-docs/stable/>
- <https://bitbucket.org/hrojas/learn-pandas>
- Python for Data Analysis (O'Reilly) by Wes McKinney (the author of pandas)

## 1.4 matplotlib

[matplotlib](#) is the primary scientific plotting library in Python. It provides:

- Publication-quality [visualizations](#) such as line charts, histograms, and scatter plots.
- Integration in pandas to make plotting much easier.
- Interactive plotting in Jupyter notebooks for quick visualizations.
  - Requires some setup. See preamble and [%matplotlib magic command](#).
- Many GUI backends, export to PDF, SVG, JPG, PNG, BMP, GIF, etc.
- Ecosystem of libraries for more advanced plotting, e.g. [Seaborn](#)

### 1.4.1 Low-level usage

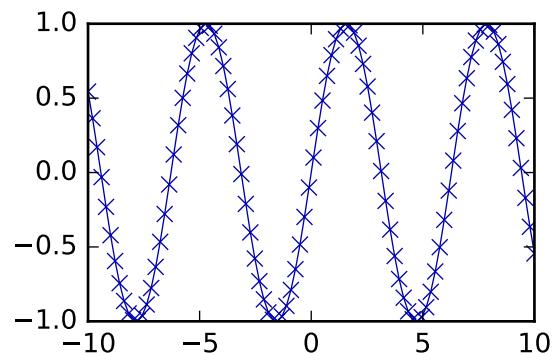
`plot()` is the [main function](#) to generate a plot (but many more exist):

```
plot(x, y)          Plot x and y using default line style and color
plot(x, y, 'bo')    Plot x and y using blue circle markers
plot(y, 'r+')       Plot y using x as index array 0..N-1, and red plusses
```

Every plotting function is completely customizable through a large set of options.

```
In [159]: plt.rcParams['figure.figsize'] = (3, 2)
          x = np.linspace(-10, 10, 100) # Sequence of integers for X-axis
          y = np.sin(x) # sine values
          plt.plot(x, y, marker="x") # Line plot with marker x
```

```
Out[159]: [<matplotlib.lines.Line2D at 0x11e8a22e8>]
```

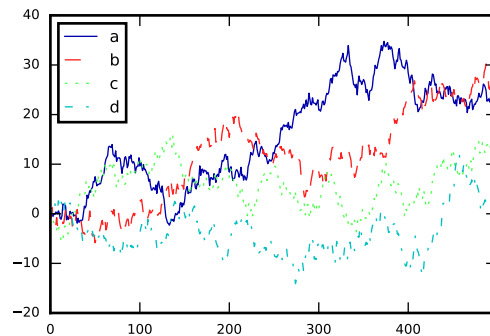


## 1.4.2 pandas + matplotlib

pandas DataFrames offer an easier, higher-level interface for matplotlib functions

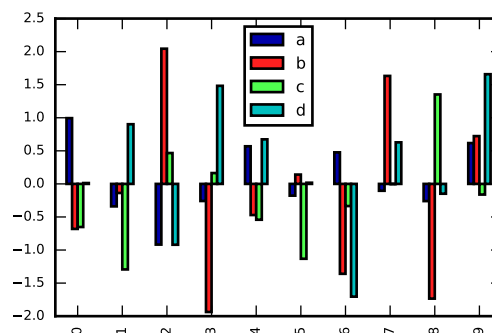
```
In [160]: plt.rcParams['font.size'] = 5; plt.rcParams['lines.linewidth'] = 0.5
          df = pd.DataFrame(np.random.randn(500, 4), columns=['a', 'b', 'c', 'd']) # random 4D data
          df.cumsum().plot() # Plot cumulative sum of all series.
```

```
Out[160]: <matplotlib.axes._subplots.AxesSubplot at 0x11e8a8518>
```



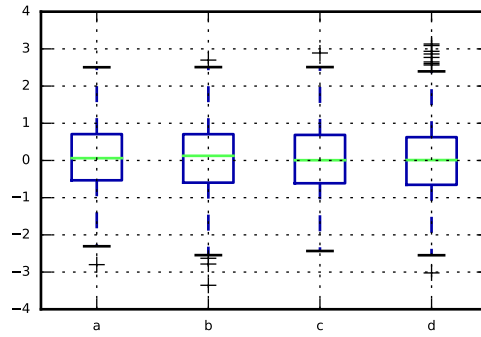
```
In [161]: df[:10].plot(kind='bar') # Plot first 10 arrays as boxplots
```

```
Out[161]: <matplotlib.axes._subplots.AxesSubplot at 0x11ebf6a20>
```



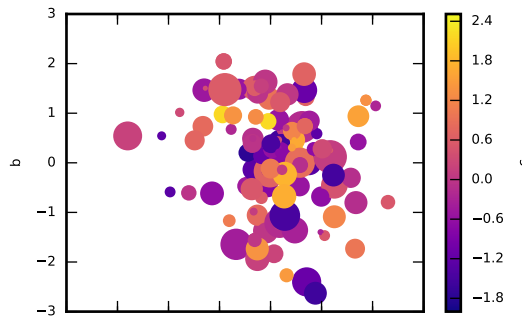
```
In [162]: df.boxplot() # Boxplot for each of the 4 series
```

```
Out[162]: <matplotlib.axes._subplots.AxesSubplot at 0x11eed9080>
```



```
In [165]: # Scatter plot using the 4 variables for x, y, color (colormap plasma), and scale, res
df[:100].plot(kind='scatter', x='a', y='b', c='c', s=df['d']*50, linewidth='0', cmap='
```

```
Out[165]: <matplotlib.axes._subplots.AxesSubplot at 0x11f7760f0>
```



### 1.4.3 Advanced plotting libraries

Several libraries, such as [Seaborn](#) offer more advanced plots and easier interfaces.

### 1.4.4 Further reading links

- [Matplotlib examples](#)
- [Plotting with pandas](#)
- [Seaborn examples](#)