

# Financial Programming and Databases

WANG, Zigan (王自干)

Ph.D., Columbia University

Faculty of Business Economics, The University of Hong Kong

# Course Outline

- Data management with Python (3 weeks)
  - Packages
    - Numpy/Scipy
    - Dataframe
    - Pandas
  - OOP
    - Class
  - Geolocation
  - Multiprocessing
  - Pickle, Json, HDF5 (.h5)
  - Web scraping (regular expression, selenium)

# External Libraries

# External libraries

A very complete list can be found at PyPi the Python Package Index:

<https://pypi.python.org/pypi>

To install, use pip, which comes with Python:

`pip install package`

or download, unzip, and run the installer directly from the directory:

`python setup.py install`

If you have Python 2 and Python 3 installed, use pip3 (though not with Anaconda) or make sure the right version is first in your PATH.

# Python Libraries for Data Science

## *NumPy:*

- introduces objects for multidimensional arrays and matrices, as well as functions that allow to easily perform advanced mathematical and statistical operations on those objects
- provides vectorization of mathematical operations on arrays and matrices which significantly improves the performance
- many other python libraries are built on NumPy

Link: <http://www.numpy.org/>

# Numpy

<http://www.numpy.org/>

Mathematics and statistics, especially multi-dimensional array manipulation for data processing.

Good introductory tutorials by Software Carpentry:

<http://swcarpentry.github.io/python-novice-inflammation/>

# Numpy data

Perhaps the nicest thing about numpy is its handling of complicated 2D datasets. It has its own array types which overload the indexing operators. Note the difference in the below from the standard [1d][2d] notation:

```
import numpy
data = numpy.int_([
[1,2,3,4,5],
[10,20,30,40,50],
[100,200,300,400,500]
])

print(data[0,0])          #      1
print(data[1:3,1:3])      #      [[20 30] [200 300]]
```

On a standard list, `data[1:3][1:3]` wouldn't work, at best `data[1:3][0][1:3]` would give you `[20][30]`

# Numpy operations

You can additionally do maths on the arrays, including matrix manipulation.

```
import numpy  
data = numpy.int_([  
[1,2,3,4,5],  
[10,20,30,40,50],  
[100,200,300,400,500]  
])  
print(data[1:3,1:3] - 10) # [[10 20], [190 290]]  
print(numpy.transpose(data[1:3,1:3])) # [[20 200], [30 300]]
```

# Python Libraries for Data Science

*SciPy:*

- collection of algorithms for linear algebra, differential equations, numerical integration, optimization, statistics and more
- part of SciPy Stack
- built on NumPy

Link: <https://www.scipy.org/scipylib/>

# Python Libraries for Data Science

## *Pandas:*

- adds data structures and tools designed to work with table-like data (similar to Series and Data Frames in R)
- provides tools for data manipulation: reshaping, merging, sorting, slicing, aggregation etc.
- allows handling missing data

Link: <http://pandas.pydata.org/>

# Pandas

<http://pandas.pydata.org/>

Data analysis.

Based on Numpy, but adds more sophistication.

# Pandas data

Pandas data focuses around DataFrames, 2D arrays with addition abilities to name and use rows and columns.

```
import pandas
df = pandas.DataFrame(
    data,
                # numpy array from before.
    index=['i','ii','iii'],
    columns=['A','B','C','D','E']
)
print (data['A'])
print(df.mean(0) ['A'])
print(df.mean(1) ['i'])
```

Prints:

```
i      1
ii     10
iii    100
Name: A, dtype: int32
37.0
3.0
```

# Python Libraries for Data Science

## *SciKit-Learn:*

- provides machine learning algorithms: classification, regression, clustering, model validation etc.
- built on NumPy, SciPy and matplotlib

Link: <http://scikit-learn.org/>

scikit-learn

<http://scikit-learn.org/>

Scientific analysis and machine learning.

Used for machine learning. Founded on Numpy data formats.

# Python Libraries for Data Science

## *matplotlib:*

- python 2D plotting library which produces publication quality figures in a variety of hardcopy formats
- a set of functionalities similar to those of MATLAB
- line plots, scatter plots, barcharts, histograms, pie charts etc.
- relatively low-level; some effort needed to create advanced visualization

Link: <https://matplotlib.org/>

# Python Libraries for Data Science

## *Seaborn:*

- based on matplotlib
- provides high level interface for drawing attractive statistical graphics
- Similar (in style) to the popular ggplot2 library in R

Link: <https://seaborn.pydata.org/>

# Beautiful Soup

<https://www.crummy.com/software/BeautifulSoup/>

Web analysis.

Need other packages to actually download pages like the library requests.

<http://docs.python-requests.org/en/master/>

BeautifulSoup navigates the Document Object Model:

<http://www.w3schools.com/>

Not a library, but a nice intro to web programming with Python.

<https://wiki.python.org/moin/WebProgramming>

# Tweepy

<http://www.tweepy.org/>

Downloading Tweets for analysis.

You'll also need a developer key:

<http://themepacific.com/how-to-generate-api-key-consumer-token-access-key-for-twitter-oauth/994/>

Most social media sites have equivalent APIs (functions to access them) and modules to use those.

NLTK

<http://www.nltk.org/>

Natural Language Toolkit.

Parse text and analyse everything from Parts Of Speech to positivity or negativity of statements (sentiment analysis).

# Celery

<http://www.celeryproject.org/>

Concurrent computing / parallelisation.

For splitting up programs and running them on multiple computers e.g. to remove memory limits.

See also:

<https://docs.python.org/3/library/concurrency.html>

# Loading Python Libraries

```
In [ ]: #Import Python Libraries  
import numpy as np  
import scipy as sp  
import pandas as pd  
import matplotlib as mpl  
import seaborn as sns
```

Press Shift+Enter to execute the *jupyter* cell

# Start Jupyter notebook

# On the Shared Computing Cluster

```
[scc1 ~] jupyter notebook
```

The screenshot shows the Jupyter Notebook web interface. At the top, there is a navigation bar with tabs for 'Files' (selected), 'Running', and 'Clusters'. On the right side of the header, there are 'Logout' and other user-related buttons. Below the header, a message says 'Select items to perform actions on them.' To the right of this message are buttons for 'Upload', 'New ▾', and a refresh icon. The main area displays a list of files:

	Name	Last Modified
<input type="checkbox"/>	<u>dataScience.ipynb</u>	8 minutes ago
<input type="checkbox"/>	flights.csv	2 minutes ago
<input type="checkbox"/>	Salaries.csv	a minute ago

# Course Outline

- Data management with Python (3 weeks)
  - Packages
    - Numpy/Scipy
    - Dataframe
    - Pandas
  - OOP
    - Class
  - Geolocation
  - Multiprocessing
  - Pickle, Json, HDF5 (.h5)
  - Web scraping (regular expression, selenium)

# NumPy

- NumPy is a Python C extension library for array-oriented computing
  - Efficient
  - In-memory
  - Contiguous (or Strided)
  - Homogeneous (but types can be algebraic)



- NumPy is suited to many applications
  - Image processing
  - Signal processing
  - Linear algebra
  - A plethora of others

# NumPy Ecosystem

OpenCV

PySAL

numexpr

astropy

PyTables

statsmodels

Biopython

scikit-image

scikit-learn

Numba

Scipy

Pandas

Matplotlib

NumPy

# Quick Start

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

In [2]: a = np.array([1,2,3,4,5,6,7,8,9])

In [3]: a
Out[3]: array([1, 2, 3, 4, 5, 6, 7, 8, 9])

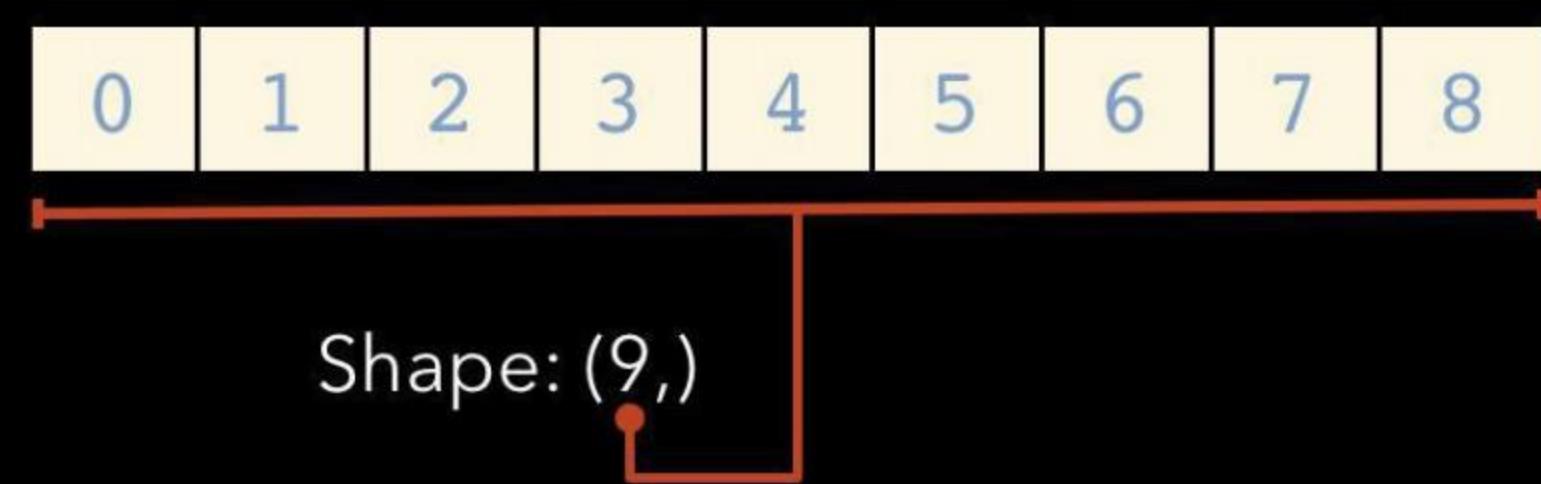
In [4]: b = a.reshape((3,3))

In [5]: b
Out[5]:
array([[1, 2, 3],
       [4, 5, 6],
       [7, 8, 9]])

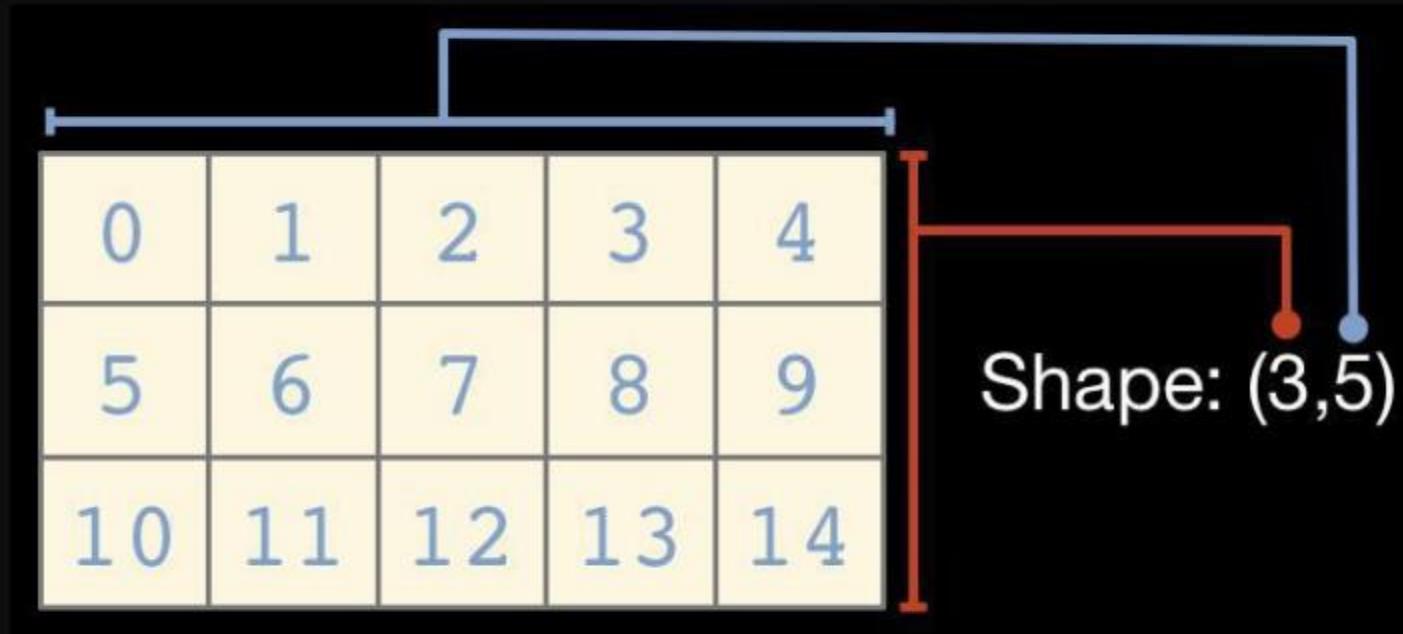
In [6]: b * 10 + 4
Out[6]:
array([[14, 24, 34],
       [44, 54, 64],
       [74, 84, 94]])
```

# Array Shape

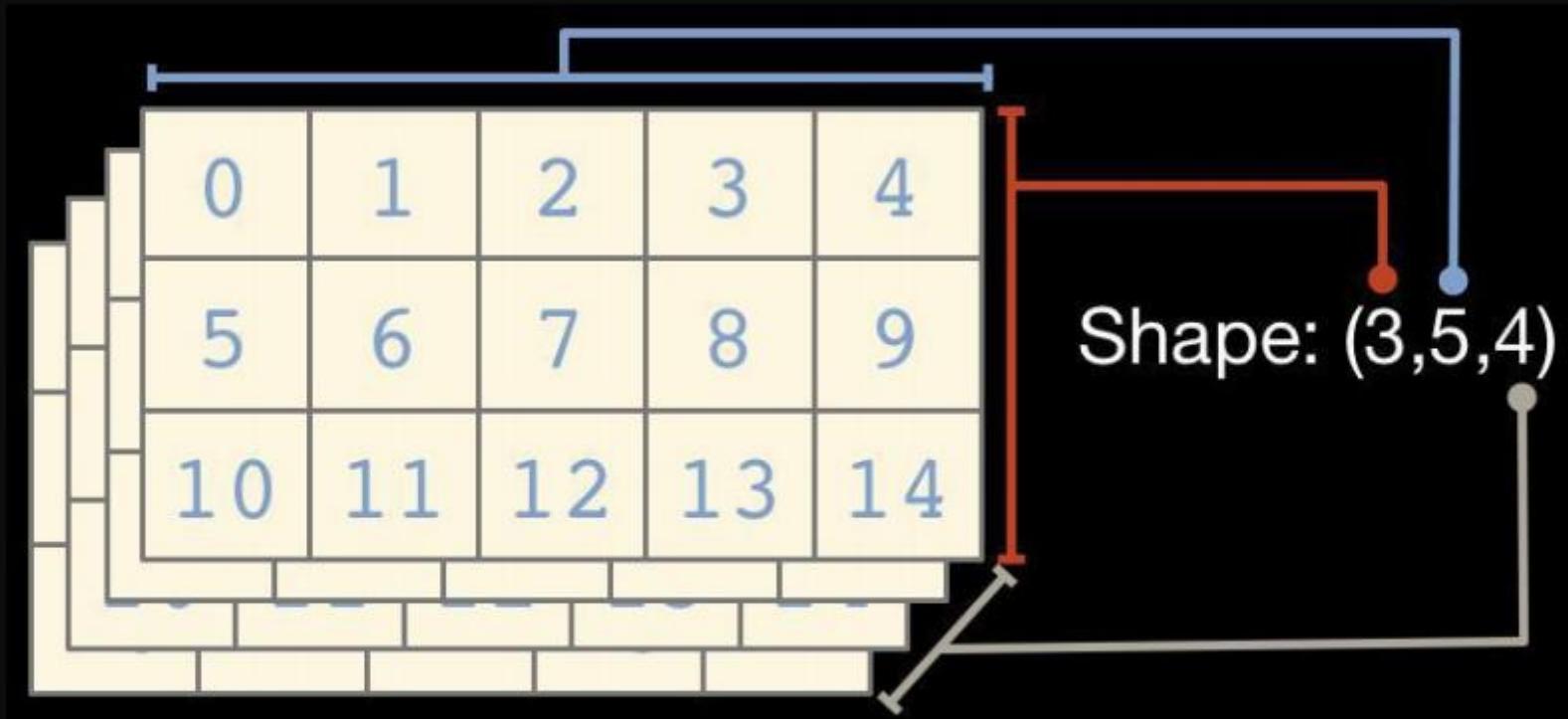
**One dimensional arrays have a 1-tuple for their shape**



**...Two dimensional arrays have a 2-tuple**



**...And so on**



# Array Creation

Explicitly from a list of values

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

As a range of values

```
In [3]: np.arange(10)
Out[3]: array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
```

By specifying the number of elements

```
In [4]: np.linspace(0, 1, 5)
Out[4]: array([ 0. ,  0.25,  0.5 ,  0.75,  1. ])
```

## Zero-initialized

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

## One-initialized

```
In [5]: np.ones((1,5))
Out[5]: array([[ 1.,  1.,  1.,  1.,  1.]])
```

## Uninitialized

```
In [4]: np.empty((1,3))
Out[4]: array([[ 2.12716633e-314,   2.12716633e-314,   2.15203762e-314]])
```

## Constant diagonal value

```
In [6]: np.eye(3)
Out[6]:
array([[ 1.,  0.,  0.],
       [ 0.,  1.,  0.],
       [ 0.,  0.,  1.]])
```

## Multiple diagonal values

```
In [7]: np.diag([1,2,3,4])
Out[7]:
array([[1, 0, 0, 0],
       [0, 2, 0, 0],
       [0, 0, 3, 0],
       [0, 0, 0, 4]])
```

# Indexing and Slicing

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

all values

arr[0:2,:]

arr[2,1:]

Implied end

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

arr[:2, 2:3]

Implied zero

NumPy array indices can also take an optional stride

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

`arr[:,::2]`

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

`arr[::2,:3]`

```
a = np.arange(10).reshape((2,5))

a.ndim      # 2 dimension
a.shape     # (2, 5) shape of array
a.size      # 10 # of elements
a.T         # transpose
a.dtype     # data type
```

## Arithmetic operators: **elementwise** application

```
a = np.arange(4)
# array([0, 1, 2, 3])

b = np.array([2, 3, 2, 4])

a * b # array([ 0,  3,  4, 12])
b - a # array([2, 2, 0, 1])

c = [2, 3, 4, 5]
a * c # array([ 0,  3,  8, 15])
```

Also, we can use `+=` and `*=`.

`a += b:`      `a = a+b`

# Introducing NumPy Arrays

## SIMPLE ARRAY CREATION

```
>>> a = array([0,1,2,3])
>>> a
array([0, 1, 2, 3])
```

## CHECKING THE TYPE

```
>>> type(a)
<type 'array'>
```

## NUMERIC ‘TYPE’ OF ELEMENTS

```
>>> a.dtype
dtype('int32')
```

## BYTES PER ELEMENT

```
>>> a.itemsize # per element
4
```

## ARRAY SHAPE

```
# shape returns a tuple
# listing the length of the
# array along each dimension.
>>> a.shape
(4,)
>>> shape(a)
(4,)
```

## ARRAY SIZE

```
# size reports the entire
# number of elements in an
# array.
>>> a.size
4
>>> size(a)
4
```

# Introducing NumPy Arrays

## BYTES OF MEMORY USED

```
# returns the number of bytes
# used by the data portion of
# the array.

>>> a.nbytes
12
```

## NUMBER OF DIMENSIONS

```
>>> a.ndim
1
```

## ARRAY COPY

```
# create a copy of the array
>>> b = a.copy()
>>> b
array([0, 1, 2, 3])
```

## CONVERSION TO LIST

```
# convert a numpy array to a
# python list.

>>> a.tolist()
[0, 1, 2, 3]
```

```
# For 1D arrays, list also
# works equivalently, but
# is slower.

>>> list(a)
[0, 1, 2, 3]
```

# Setting Array Elements

## ARRAY INDEXING

```
>>> a[0]
0
>>> a[0] = 10
>>> a
[10, 1, 2, 3]
```

## FILL

```
# set all values in an array.
>>> a.fill(0)
>>> a
[0, 0, 0, 0]

# This also works, but may
# be slower.
>>> a[:] = 1
>>> a
[1, 1, 1, 1]
```



## BEWARE OF TYPE COERSION

```
>>> a.dtype
dtype('int32')

# assigning a float to into
# an int32 array will
# truncate decimal part.
>>> a[0] = 10.6
>>> a
[10, 1, 2, 3]

# fill has the same behavior
>>> a.fill(-4.8)
>>> a
[-4, -4, -4, -4]
```

# Multi-Dimensional Arrays

## MULTI-DIMENSIONAL ARRAYS

```
>>> a = array([[ 0,  1,  2,  3],
   [10,11,12,13]])  

>>> a  

array([[ 0,  1,  2,  3],
       [10,11,12,13]])
```

## (ROWS,COLUMNS)

```
>>> a.shape  

(2, 4)  

>>> shape(a)  

(2, 4)
```

## ELEMENT COUNT

```
>>> a.size  

8  

>>> size(a)  

8
```

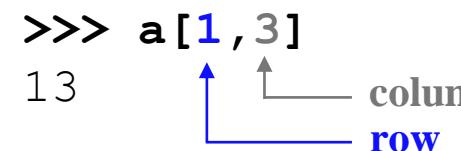
## NUMBER OF DIMENSIONS

```
>>> a.ndim  

2
```

## GET/SET ELEMENTS

```
>>> a[1,3]
```



row  
column

```
>>> a[1,3] = -1  

>>> a  

array([[ 0,  1,  2,  3],
       [10,11,12, -1]])
```

## ADDRESS FIRST ROW USING SINGLE INDEX

```
>>> a[1]  

array([10, 11, 12, -1])
```

## Array broadcasting with scalars

This also allows us to add a constant to a matrix or multiply a matrix by a constant

```
A = np.ones((3,3))

print 3 * A - 1
# [[ 2.  2.  2.]
#  [ 2.  2.  2.]
#  [ 2.  2.  2.]]
```

# Broadcasting

A key feature of NumPy is broadcasting, where arrays with different, but compatible shapes can be used as arguments to ufuncs

a	0	1	2	3	4
	10	10	10	10	10
c	10	11	12	13	14

$c = a + 10$

In this case an array scalar is broadcast to an array with shape (5, )

A slightly more involved broadcasting example in two dimensions

$$c = a + b$$

0	1		
2	3	+	
4	5		=
a	b		c

10	10
20	20
30	30

0	11
22	23
34	35

Here an array of shape  $(3, 1)$  is broadcast to an array with shape  $(3, 2)$

# Array broadcasting

$$\begin{array}{c} \begin{array}{ccc} 0 & 0 & 0 \\ 10 & 10 & 10 \\ 20 & 20 & 20 \\ 30 & 30 & 30 \end{array} + \begin{array}{ccc} 0 & 1 & 2 \\ 0 & 1 & 2 \\ 0 & 1 & 2 \\ 0 & 1 & 2 \end{array} = \begin{array}{ccc} 0 & 0 & 0 \\ 10 & 10 & 10 \\ 20 & 20 & 20 \\ 30 & 30 & 30 \end{array} + \begin{array}{ccc} 0 & 1 & 2 \\ 0 & 1 & 2 \\ 0 & 1 & 2 \\ 0 & 1 & 2 \end{array} = \end{array} \quad \begin{array}{c} \nearrow \\ \begin{array}{ccc} 0 & 0 & 0 \\ 10 & 10 & 10 \\ 20 & 20 & 20 \\ 30 & 30 & 30 \end{array} + \begin{array}{c} 0 \ 1 \ 2 \end{array} = \begin{array}{ccc} 0 & 0 & 0 \\ 10 & 10 & 10 \\ 20 & 20 & 20 \\ 30 & 30 & 30 \end{array} + \begin{array}{ccc} 0 & 1 & 2 \\ 0 & 1 & 2 \\ 0 & 1 & 2 \\ 0 & 1 & 2 \end{array} = \begin{array}{ccc} 0 & 1 & 2 \\ 10 & 11 & 12 \\ 20 & 21 & 22 \\ 30 & 31 & 32 \end{array} \end{array} \\ \quad \begin{array}{c} \uparrow \\ \begin{array}{c} 0 \\ 10 \\ 20 \\ 30 \end{array} + \begin{array}{c} 0 \ 1 \ 2 \end{array} = \begin{array}{ccc} 0 & 0 & 0 \\ 10 & 10 & 10 \\ 20 & 20 & 20 \\ 30 & 30 & 30 \end{array} + \begin{array}{ccc} 0 & 1 & 2 \\ 0 & 1 & 2 \\ 0 & 1 & 2 \\ 0 & 1 & 2 \end{array} = \end{array} \end{array} \end{array}$$

# Broadcasting Rules

In order for an operation to broadcast, the size of all the trailing dimensions for both arrays must either:

be **equal** OR be **one**

A (1d array): 3

B (2d array): 2 x 3

Result (2d array): 2 x 3

A (2d array): 6 x 1

B (3d array): 1 x 6 x 4

Result (3d array): 1 x 6 x 4

A (4d array): 3 x 1 x 6 x 1

B (3d array): 2 x 1 x 4

Result (4d array): 3 x 2 x 6 x 4

## Vector operations

- inner product
- outer product
- dot product (matrix multiplication)

```
# note: numpy automatically converts lists
u = [1, 2, 3]
v = [1, 1, 1]

np.inner(u, v)
# 6
np.outer(u, v)
# array([[1, 1, 1],
#        [2, 2, 2],
#        [3, 3, 3]])
np.dot(u, v)
# 6
```

## Matrix operations

First, define some matrices:

```
A = np.ones((3, 2))
# array([[ 1.,  1.],
#        [ 1.,  1.],
#        [ 1.,  1.]])
A.T
# array([[ 1.,  1.,  1.],
#        [ 1.,  1.,  1.]])  
  
B = np.ones((2, 3))
# array([[ 1.,  1.,  1.],
#        [ 1.,  1.,  1.]])
```

## Matrix operations

```
np.dot(A, B)
# array([[ 2.,  2.,  2.],
#        [ 2.,  2.,  2.],
#        [ 2.,  2.,  2.]])  
  
np.dot(B, A)
# array([[ 3.,  3.],
#        [ 3.,  3.]])  
  
np.dot(B.T, A.T)
# array([[ 2.,  2.,  2.],
#        [ 2.,  2.,  2.],
#        [ 2.,  2.,  2.]])  
  
np.dot(A, B.T)
# Traceback (most recent call last):
#   File "<stdin>", line 1, in <module>
# ValueError: shapes (3,2) and (3,2) not aligned: ...  
# ... 2 (dim 1) != 3 (dim 0)
```

# Universal Functions (ufuncs)

NumPy ufuncs are functions that operate element-wise on one or more arrays

a	0	1	2	3	4
---	---	---	---	---	---

b	0	10	20	30	40
---	---	----	----	----	----

$$c = a + b$$

c	0	11	22	33	44
---	---	----	----	----	----

ufuncs dispatch to optimized C inner-loops based on array dtype

# Mathematic Binary Operators

<code>a + b → add(a,b)</code>
<code>a - b → subtract(a,b)</code>
<code>a % b → remainder(a,b)</code>

<code>a * b → multiply(a,b)</code>
<code>a / b → divide(a,b)</code>
<code>a ** b → power(a,b)</code>

## MULTIPLY BY A SCALAR

```
>>> a = array((1,2))
>>> a*3.
array([3.,  6.])
```

## ELEMENT BY ELEMENT ADDITION

```
>>> a = array([1,2])
>>> b = array([3,4])
>>> a + b
array([4, 6])
```

## ADDITION USING AN OPERATOR FUNCTION

```
>>> add(a,b)
array([4, 6])
```

## ⚠ IN PLACE OPERATION

```
# Overwrite contents of a.
# Saves array creation
# overhead
>>> add(a,b,a) # a += b
array([4, 6])
>>> a
array([4, 6])
```

# Comparison and Logical Operators

`equal`               `(==)`  
`greater_equal` `(>=)`  
`logical_and`  
`logical_not`

`not_equal`     `(!=)`  
`less`              `(<)`  
`logical_or`

`greater`          `(>)`  
`less_equal`     `(<=)`  
`logical_xor`

## 2D EXAMPLE

```
>>> a = array(((1,2,3,4),(2,3,4,5)))  
>>> b = array(((1,2,5,4),(1,3,4,5)))  
>>> a == b  
array([[True, True, False, True],  
       [False, True, True, True]])  
# functional equivalent  
>>> equal(a,b)  
array([[True, True, False, True],  
       [False, True, True, True]])
```

# Bitwise Operators

<code>bitwise_and (&amp;)</code>	<code>invert (~)</code>	<code>right_shift(a,shifts)</code>
<code>bitwise_or ( )</code>	<code>bitwise_xor</code>	<code>left_shift (a,shifts)</code>

## BITWISE EXAMPLES

```
>>> a = array((1,2,4,8))
>>> b = array((16,32,64,128))
>>> bitwise_or(a,b)
array([ 17,   34,   68,  136])

# bit inversion
>>> a = array((1,2,3,4), uint8)
>>> invert(a)
array([254, 253, 252, 251], dtype=uint8)

# left shift operation
>>> left_shift(a,3)
array([ 8, 16, 24, 32], dtype=uint8)
```

# Trig and Other Functions

## TRIGONOMETRIC

<code>sin(x)</code>	<code>sinh(x)</code>
<code>cos(x)</code>	<code>cosh(x)</code>
<code>arccos(x)</code>	<code>arccosh(x)</code>
<code>arctan(x)</code>	<code>arctanh(x)</code>
<code>arcsin(x)</code>	<code>arcsinh(x)</code>
<code>arctan2(x,y)</code>	

## OTHERS

<code>exp(x)</code>	<code>log(x)</code>
<code>log10(x)</code>	<code>sqrt(x)</code>
<code>absolute(x)</code>	<code>conjugate(x)</code>
<code>negative(x)</code>	<code>ceil(x)</code>
<code>floor(x)</code>	<code>fabs(x)</code>
<code>hypot(x,y)</code>	<code>fmod(x,y)</code>
<code>maximum(x,y)</code>	<code>minimum(x,y)</code>

## hypot(x,y)

Element by element distance  
calculation using  $\sqrt{x^2 + y^2}$

# NumPy has many built-in ufuncs

- comparison: `<`, `<=`, `==`, `!=`, `>=`, `>`
- arithmetic: `+`, `-`, `*`, `/`, `reciprocal`, `square`
- exponential: `exp`, `expm1`, `exp2`, `log`, `log10`, `log1p`, `log2`,  
`power`, `sqrt`
- trigonometric: `sin`, `cos`, `tan`, `acsin`, `arccos`, `atctan`
- hyperbolic: `sinh`, `cosh`, `tanh`, `acsinh`, `arccosh`, `atctanh`
- bitwise operations: `&`, `|`, `~`, `^`, `left_shift`, `right_shift`
- logical operations: `and`, `logical_xor`, `not`, `or`
- predicates: `isfinite`, `isinf`, `isnan`, `signbit`
- other: `abs`, `ceil`, `floor`, `mod`, `modf`, `round`, `sinc`, `sign`,  
`trunc`

# Axis

Array method reductions take an optional `axis` parameter that specifies over which axes to reduce

`axis=None` reduces into a single scalar

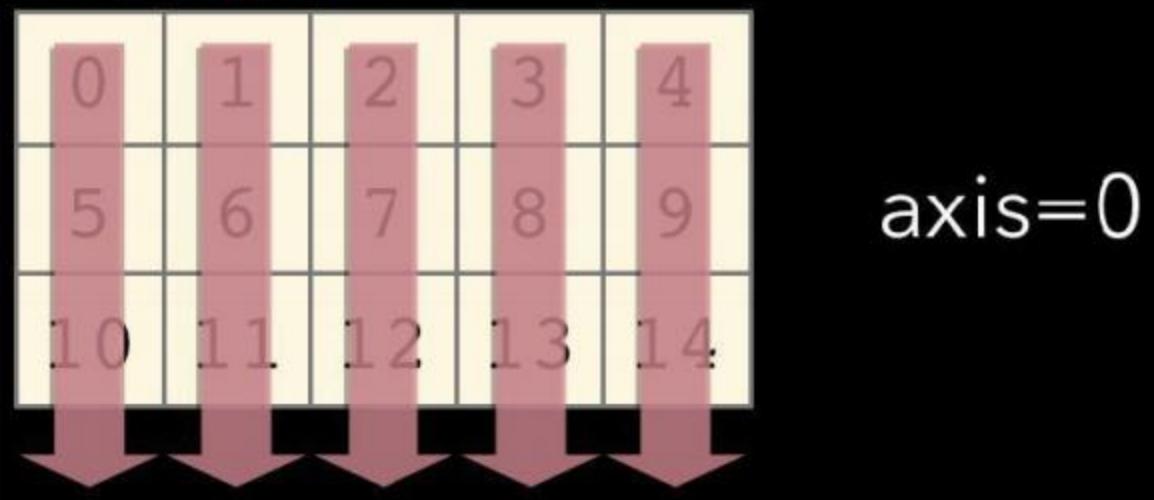
```
In [7]: a.sum  
()  
Out[7]: 105
```

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

`axis=None`

`axis=0` reduces into the zeroth dimension

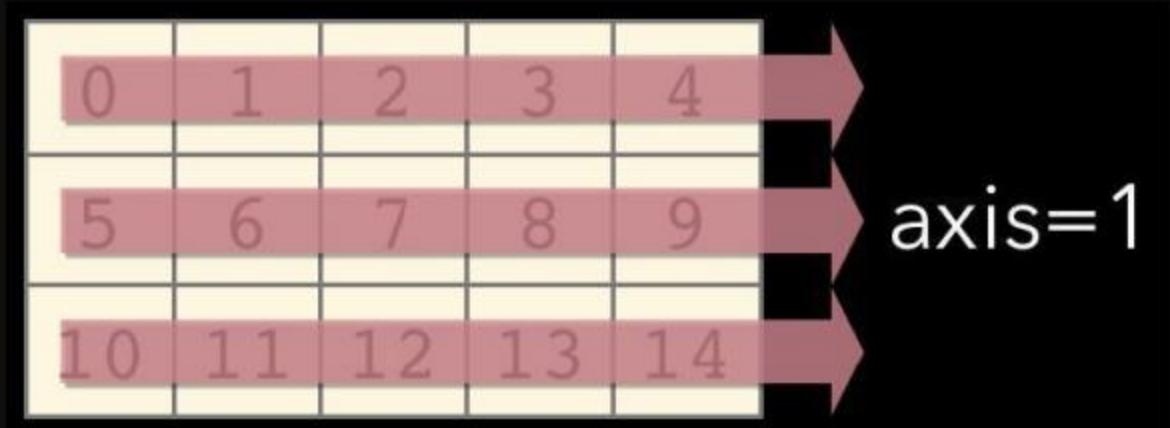
```
In [8]: a.sum(axis=0)  
Out[8]: array([15, 18, 21, 24,  
27])
```



`axis=0`

`axis=1` reduces into the first dimension

```
In [9]: a.sum(axis=1)  
Out[9]: array([10, 35, 60])
```



`axis=1`

# Fancy Indexing

NumPy arrays may be used to index into other arrays

```
In [2]: a = np.arange(15).reshape((3,5))

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

In [5]: j = np.array([[2, 1], [4, 4]])

In [6]: a[i,j]
Out[6]:
array([[ 2,  6],
       [ 9, 14]])
```

Boolean arrays can also be used as indices into other arrays

```
In [2]: a = np.arange(15).reshape((3,5))

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

In [4]: b = (a % 3 == 0)

In [5]: b
Out[5]:
array([[ True, False, False,  True, False],
       [False,  True, False, False,  True],
       [False, False,  True, False, False]], dtype=bool)

In [6]: a[b]
Out[6]: array([ 0,  3,  6,  9, 12])
```

# NumPy Functions

- Data I/O
  - `fromfile`, `genfromtxt`, `load`, `loadtxt`, `save`, `savetxt`
- Mesh Creation
  - `mgrid`, `meshgrid`, `ogrid`
- Manipulation
  - `einsum`, `hstack`, `take`, `vstack`

- `numpy.fft` — Fast Fourier transforms
- `numpy.polynomial` — Efficient polynomials
- `numpy.linalg` — Linear algebra
  - `cholesky`, `det`, `eig`, `eigvals`, `inv`, `lstsq`, `norm`, `qr`, `svd`
- `numpy.math` — C standard library math functions
- `numpy.random` — Random number generation
  - `beta`, `gamma`, `geometric`, `hypergeometric`, `lognormal`, `normal`, `poisson`, `uniform`, `weibull`

loc: mean  
scale: std. dev

```
np.random.random((2,3))
# array([[ 0.78084261,  0.64328818,  0.55380341],
#        [ 0.24611092,  0.37011213,  0.83313416]])

a = np.random.normal(loc=1.0, scale=2.0, size=(2,2))
# array([[ 2.87799514,  0.6284259 ],
#        [ 3.10683164,  2.05324587]])

np.savetxt("a_out.txt", a)
# save to file
b = np.loadtxt("a_out.txt")
# read from file
```

## Operations along axes

```
a = np.random.random((2,3))
# array([[ 0.9190687 ,  0.36497813,  0.75644216],
#        [ 0.91938241,  0.08599547,  0.49544003]])
a.sum()
# 3.5413068994445549
a.sum(axis=0) # column sum
# array([ 1.83845111,  0.4509736 ,  1.25188219])
a.cumsum()
# array([ 0.9190687 ,  1.28404683,  2.04048899,  2.9598714 ,
#        3.04586687,  3.5413069 ])
a.cumsum(axis=1) # cumulative row sum
# array([[ 0.9190687 ,  1.28404683,  2.04048899],
#        [ 0.91938241,  1.00537788,  1.50081791]])
a.min()
# 0.0859954690403677
a.max(axis=0)
# array([ 0.91938241,  0.36497813,  0.75644216])
```

# Array Calculation Methods

## SUM FUNCTION

```
>>> a = array([[1,2,3],  
              [4,5,6]], float)  
  
# Sum defaults to summing all  
# *all* array values.  
>>> sum(a)  
21.  
  
# supply the keyword axis to  
# sum along the 0th axis.  
>>> sum(a, axis=0)  
array([5., 7., 9.])  
  
# supply the keyword axis to  
# sum along the last axis.  
>>> sum(a, axis=-1)  
array([6., 15.])
```

## SUM ARRAY METHOD

```
# The a.sum() defaults to  
# summing *all* array values
```

```
>>> a.sum()  
21.
```

```
# Supply an axis argument to  
# sum along a specific axis.
```

```
>>> a.sum(axis=0)  
array([5., 7., 9.])
```

## PRODUCT

```
# product along columns.
```

```
>>> a.prod(axis=0)  
array([ 4., 10., 18.])
```

```
# functional form.
```

```
>>> prod(a, axis=0)  
array([ 4., 10., 18.])
```

# Min/Max

## MIN

```
>>> a = array([2.,3.,0.,1.])
>>> a.min(axis=0)
0.
# use Numpy's amin() instead
# of Python's builtin min()
# for speed operations on
# multi-dimensional arrays.
>>> amin(a, axis=0)
0.
```

## ARGMIN

```
# Find index of minimum value.
>>> a.argmin(axis=0)
2
# functional form
>>> argmin(a, axis=0)
2
```

## MAX

```
>>> a = array([2.,1.,0.,3.])
>>> a.max(axis=0)
3.
# functional form
>>> amax(a, axis=0)
3.
```

## ARGMAX

```
# Find index of maximum value.
>>> a.argmax(axis=0)
1
# functional form
>>> argmax(a, axis=0)
1
```

# Statistics Array Methods

## MEAN

```
>>> a = array([[1,2,3],  
             [4,5,6]], float)  
  
# mean value of each column  
>>> a.mean(axis=0)  
array([ 2.5,  3.5,  4.5])  
>>> mean(a, axis=0)  
array([ 2.5,  3.5,  4.5])  
>>> average(a, axis=0)  
array([ 2.5,  3.5,  4.5])  
  
# average can also calculate  
# a weighted average  
>>> average(a, weights=[1,2],  
...           axis=0)  
array([ 3.,  4.,  5.])
```

## STANDARD DEV./VARIANCE

```
# Standard Deviation  
>>> a.std(axis=0)  
array([ 1.5,  1.5,  1.5])  
  
# Variance  
>>> a.var(axis=0)  
array([2.25, 2.25, 2.25])  
>>> var(a, axis=0)  
array([2.25, 2.25, 2.25])
```

# Other Array Methods

## CLIP

```
# Limit values to a range
>>> a = array([[1,2,3],
              [4,5,6]], float)

# Set values < 3 equal to 3.
# Set values > 5 equal to 5.
>>> a.clip(3,5)
>>> a
array([[ 3.,  3.,  3.],
       [ 4.,  5.,  5.]])
```

## ROUND

```
# Round values in an array.
# Numpy rounds to even, so
# 1.5 and 2.5 both round to 2.
>>> a = array([1.35, 2.5, 1.5])
>>> a.round()
array([ 1.,  2.,  2.])

# Round to first decimal place.
>>> a.round(decimals=1)
array([ 1.4,  2.5,  1.5])
```

## POINT TO POINT

```
# Calculate max - min for
# array along columns
>>> a.ptp(axis=0)
array([ 3.0,  3.0,  3.0])
# max - min for entire array.
>>> a.ptp(axis=None)
5.0
```

# Summary of (most) array attributes/methods

## BASIC ATTRIBUTES

`a.dtype` – Numerical type of array elements. `float32`, `uint8`, etc.  
`a.shape` – Shape of the array. `(m,n,o,...)`  
`a.size` – Number of elements in entire array.  
`a.itemsize` – Number of bytes used by a single element in the array.  
`a.nbytes` – Number of bytes used by entire array (data only).  
`a.ndim` – Number of dimensions in the array.

## SHAPE OPERATIONS

`a.flat` – An iterator to step through array as if it is 1D.  
`a.flatten()` – Returns a 1D copy of a multi-dimensional array.  
`a.ravel()` – Same as `flatten()`, but returns a 'view' if possible.  
`a.resize(new_size)` – Change the size/shape of an array in-place.  
`a.swapaxes(axis1, axis2)` – Swap the order of two axes in an array.  
`a.transpose(*axes)` – Swap the order of any number of array axes.  
`a.T` – Shorthand for `a.transpose()`  
`a.squeeze()` – Remove any length=1 dimensions from an array.

# Summary of (most) array attributes/methods

## FILL AND COPY

`a.copy()` – Return a copy of the array.  
`a.fill(value)` – Fill array with a scalar value.

## CONVERSION / COERSION

`a.tolist()` – Convert array into nested lists of values.  
`a.tostring()` – raw copy of array memory into a python string.  
`a.astype(dtype)` – Return array coerced to given dtype.  
`a.byteswap(False)` – Convert byte order (big <-> little endian).

## COMPLEX NUMBERS

`a.real` – Return the real part of the array.  
`a.imag` – Return the imaginary part of the array.  
`a.conjugate()` – Return the complex conjugate of the array.  
`a.conj()` – Return the complex conjugate of an array. (same as conjugate)

# Summary of (most) array attributes/methods

## SAVING

```
a.dump(file) - Store a binary array data out to the given file.  
a.dumps() - returns the binary pickle of the array as a string.  
a.tofile(fid, sep="", format="%s") Formatted ascii output to file.
```

## SEARCH / SORT

```
a.nonzero() - Return indices for all non-zero elements in a.  
a.sort(axis=-1) - Inplace sort of array elements along axis.  
a.argsort(axis=-1) - Return indices for element sort order along axis.  
a.searchsorted(b) - Return index where elements from b would go in a.
```

## ELEMENT MATH OPERATIONS

```
a.clip(low, high) - Limit values in array to the specified range.  
a.round(decimals=0) - Round to the specified number of digits.  
a.cumsum(axis=None) - Cumulative sum of elements along axis.  
a.cumprod(axis=None) - Cumulative product of elements along axis.
```

# Summary of (most) array attributes/methods

## REDUCTION METHODS

All the following methods “reduce” the size of the array by 1 dimension by carrying out an operation along the specified axis. If axis is None, the operation is carried out across the entire array.

- a.sum(axis=None) – Sum up values along axis.
- a.prod(axis=None) – Find the product of all values along axis.
- a.min(axis=None) – Find the minimum value along axis.
- a.max(axis=None) – Find the maximum value along axis.
- a.argmin(axis=None) – Find the index of the minimum value along axis.
- a.argmax(axis=None) – Find the index of the maximum value along axis.
- a.ptp(axis=None) – Calculate a.max(axis) – a.min(axis)
- a.mean(axis=None) – Find the mean (average) value along axis.
- a.std(axis=None) – Find the standard deviation along axis.
- a.var(axis=None) – Find the variance along axis.
  
- a.any(axis=None) – True if any value along axis is non-zero. (or)
- a.all(axis=None) – True if all values along axis are non-zero. (and)

# Array Operations

## SIMPLE ARRAY MATH

```
>>> a = array([1,2,3,4])
>>> b = array([2,3,4,5])
>>> a + b
array([3, 5, 7, 9])
```



Numpy defines the following constants:

pi = 3.14159265359  
e = 2.71828182846

## MATH FUNCTIONS

```
# Create array from 0 to 10
>>> x = arange(11.)

# multiply entire array by
# scalar value
>>> a = (2*pi)/10.
>>> a
0.62831853071795862
>>> a*x
array([ 0., 0.628,..., 6.283])
```

```
# inplace operations
```

```
>>> x *= a
>>> x
array([ 0., 0.628,..., 6.283])
```

```
# apply functions to array.
>>> y = sin(x)
```

# Recommendations

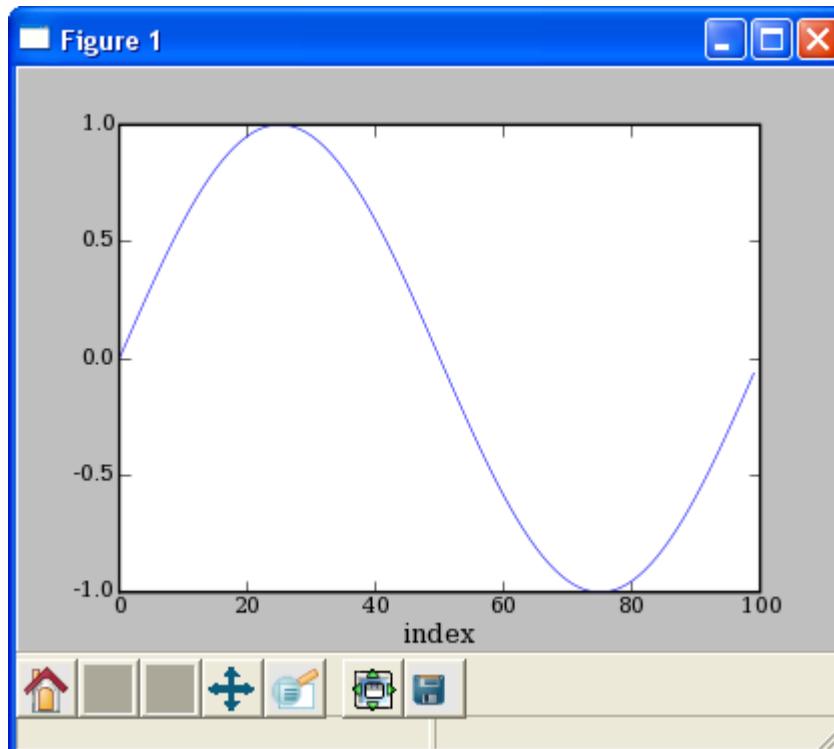
- **Matplotlib** for day-to-day data exploration.

Matplotlib has a large community, tons of plot types, and is well integrated into ipython. It is the de-facto standard for ‘command line’ plotting from ipython.

# Line Plots

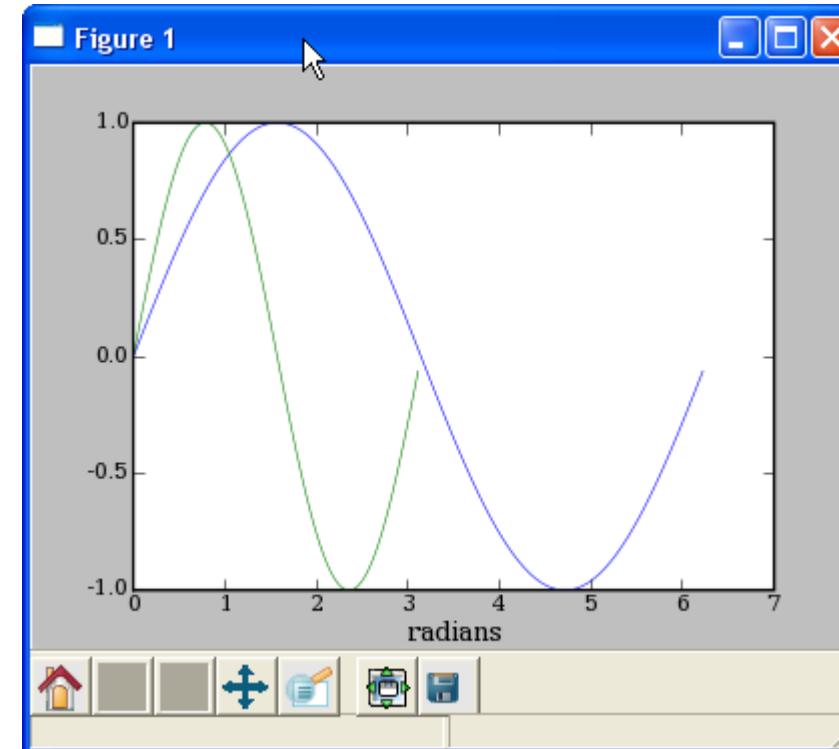
## PLOT AGAINST INDICES

```
>>> x = arange(50)*2*pi/50.
>>> y = sin(x)
>>> mpl.pyplot.plot(y)
>>> mpl.pyplot.xlabel('index')
```



## MULTIPLE DATA SETS

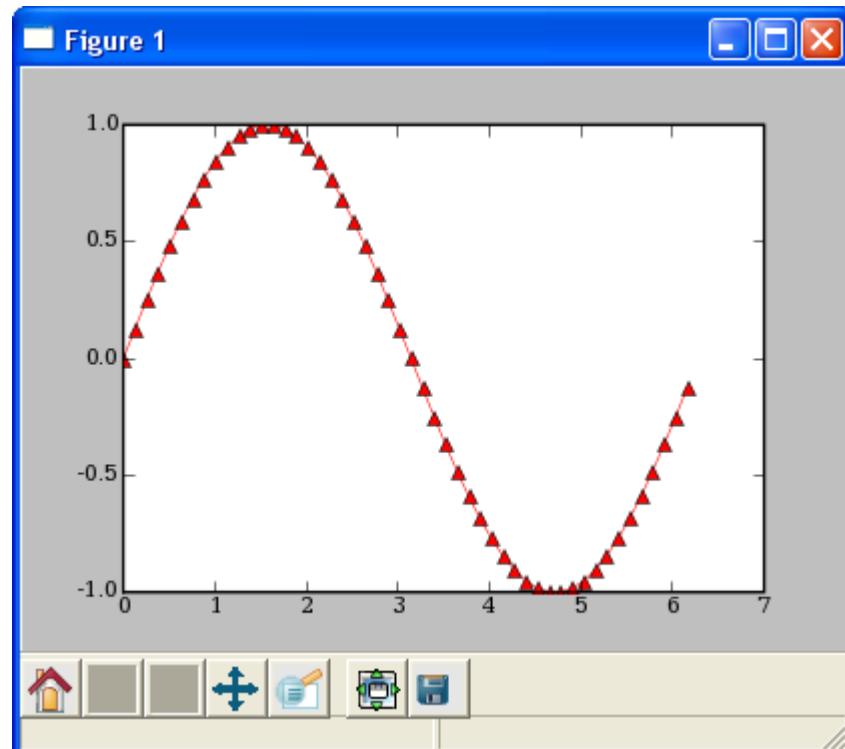
```
>>> plot(x,y,x2,y2)
>>> mpl.pyplot.xlabel('radians')
```



# Line Plots

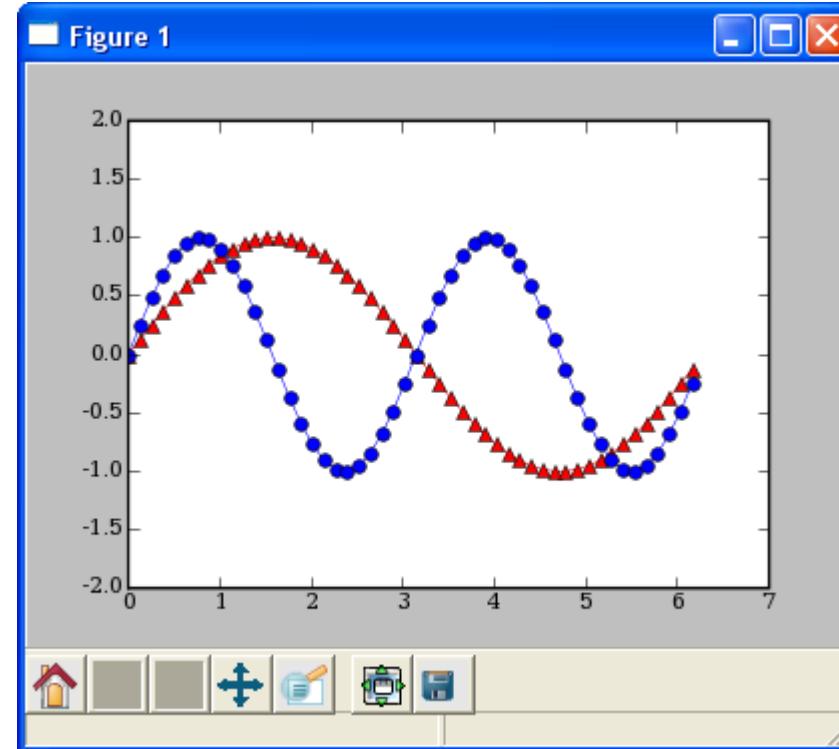
## LINE FORMATTING

```
# red, dot-dash, triangles  
>>> plot(x,sin(x), 'r-^')
```



## MULTIPLE PLOT GROUPS

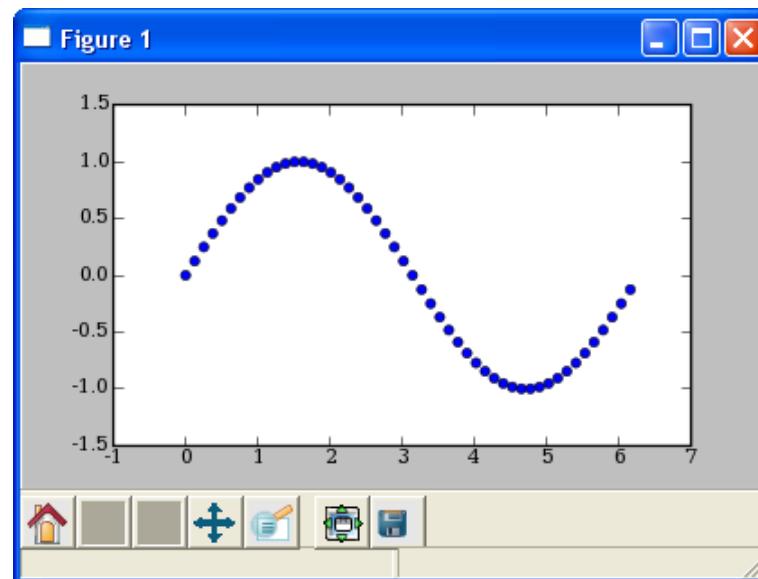
```
>>> plot(x,y1,'b-o', x,y2), r-^')  
>>> axis([0,7,-2,2])
```



# Scatter Plots

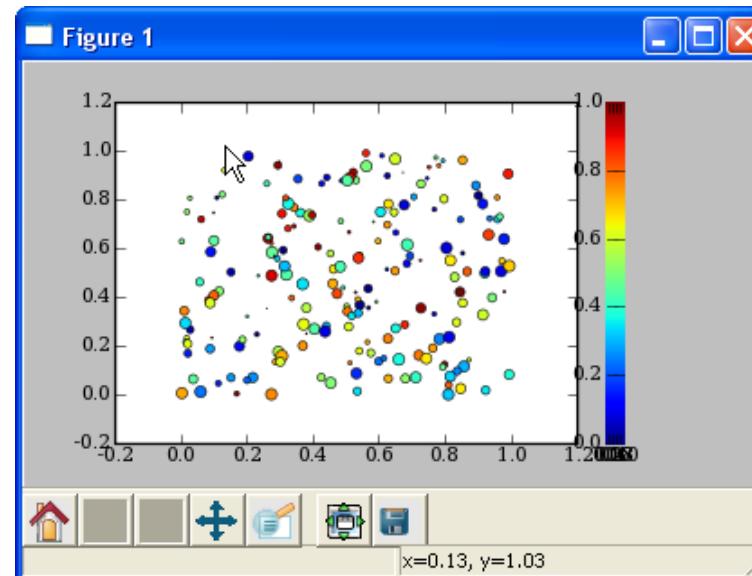
## SIMPLE SCATTER PLOT

```
>>> x = arange(50)*2*pi/50.
>>> y = sin(x)
>>> pyplot.scatter(x,y)
```



## COLORMAPPED SCATTER

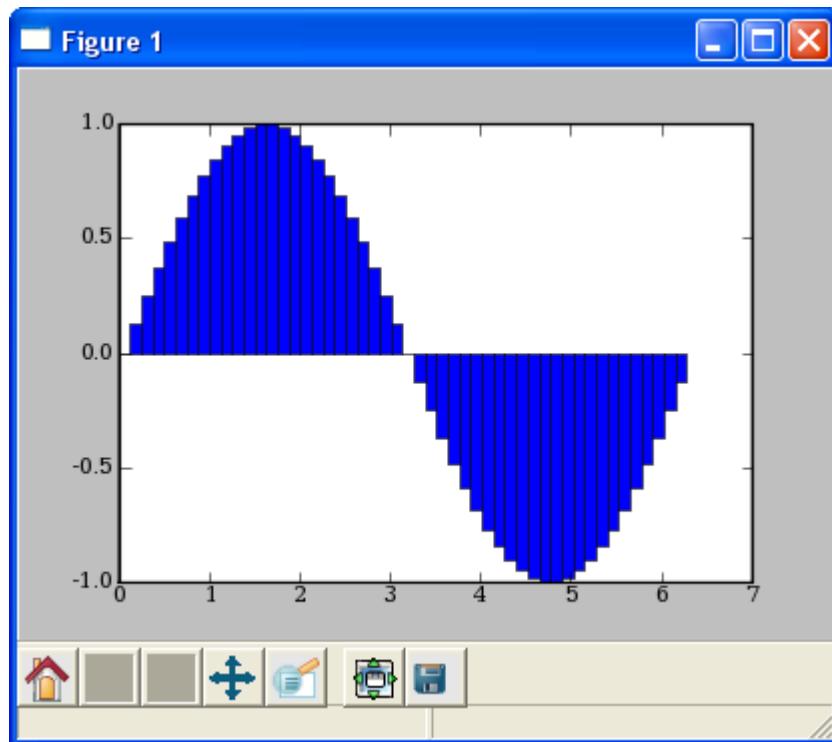
```
# marker size/color set with data
>>> x = rand(200)
>>> y = rand(200)
>>> size = rand(200)*30
>>> color = rand(200)
>>> mpl.pyplot.scatter(x, y, size,
color)
>>> colorbar()
```



# Bar Plots

## BAR PLOT

```
>>> pyplot.bar(x,sin(x),  
...           width=x[1]-x[0])
```



SciPy

# SciPy

- As the functions are less used these days, I leave this part for you to explore by yourself.

# Course Outline

- Data management with Python (3 weeks)
  - Packages
    - Numpy/Scipy
    - **Dataframe**
    - **Pandas**
  - OOP
    - Class
  - Geolocation
  - Multiprocessing
  - Pickle, Json, HDF5 (.h5)
  - Web scraping (regular expression, selenium)

# Reading data using pandas

```
In [ ]: #Read csv file  
df = pd.read_csv("http://rcs.bu.edu/examples/python/data_analysis/Salaries.csv")
```

**Note:** The above command has many optional arguments to fine-tune the data import process.

There is a number of pandas commands to read other data formats:

```
pd.read_excel('myfile.xlsx', sheet_name='Sheet1', index_col=None, na_values=['NA'])  
pd.read_stata('myfile.dta')  
pd.read_sas('myfile.sas7bdat')  
pd.read_hdf('myfile.h5', 'df')
```

# Exploring data frames

```
In [3]: #List first 5 records  
df.head()
```

Out [3] :

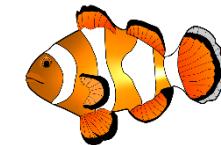
	rank	discipline	phd	service	sex	salary
0	Prof	B	56	49	Male	186960
1	Prof	A	12	6	Male	93000
2	Prof	A	23	20	Male	110515
3	Prof	A	40	31	Male	131205
4	Prof	B	20	18	Male	104800



## Hands-on exercises

- ✓ Try to read the first 10, 20, 50 records;
- ✓ Can you guess how to view the last few records;

*Hint:*



# Data Frame data types

Pandas Type	Native Python Type	Description
object	string	The most general dtype. Will be assigned to your column if column has mixed types (numbers and strings).
int64	int	Numeric characters. 64 refers to the memory allocated to hold this character.
float64	float	Numeric characters with decimals. If a column contains numbers and NaNs(see below), pandas will default to float64, in case your missing value has a decimal.
datetime64, timedelta[ns]	N/A (but see the <a href="#">datetime</a> module in Python's standard library)	Values meant to hold time data. Look into these for time series experiments.

# Data Frame data types

```
In [4]: #Check a particular column type  
df['salary'].dtype
```

```
Out[4]: dtype('int64')
```

```
In [5]: #Check types for all the columns  
df.dtypes
```

```
Out[4]: rank          object  
discipline      object  
phd            int64  
service          int64  
sex            object  
salary          int64  
dtype: object
```

# Data Frames attributes

Python objects have *attributes* and *methods*.

df.attribute	description
dtypes	list the types of the columns
columns	list the column names
axes	list the row labels and column names
ndim	number of dimensions
size	number of elements
shape	return a tuple representing the dimensionality
values	numpy representation of the data



## Hands-on exercises

- ✓ Find how many records this data frame has;
- ✓ How many elements are there?
- ✓ What are the column names?
- ✓ What types of columns we have in this data frame?

# Data Frames methods

Unlike attributes, python methods have *parenthesis*.

All attributes and methods can be listed with a `dir()` function: `dir(df)`

<b>df.method()</b>	<b>description</b>
<code>head( [n] ), tail( [n] )</code>	first/last n rows
<code>describe()</code>	generate descriptive statistics (for numeric columns only)
<code>max(), min()</code>	return max/min values for all numeric columns
<code>mean(), median()</code>	return mean/median values for all numeric columns
<code>std()</code>	standard deviation
<code>sample([n])</code>	returns a random sample of the data frame
<code>dropna()</code>	drop all the records with missing values



## Hands-on exercises

- ✓ Give the summary for the numeric columns in the dataset
- ✓ Calculate standard deviation for all numeric columns;
- ✓ What are the mean values of the first 50 records in the dataset? *Hint:* use `head()` method to subset the first 50 records and then calculate the mean

# Selecting a column in a Data Frame

*Method 1:* Subset the data frame using column name:

```
df['sex']
```

*Method 2:* Use the column name as an attribute:

```
df.sex
```

*Note:* there is an attribute *rank* for pandas data frames, so to select a column with a name "rank" we should use method 1.



## Hands-on exercises

- ✓ Calculate the basic statistics for the *salary* column;
- ✓ Find how many values in the *salary* column (use *count* method);
- ✓ Calculate the average salary;

groupby

# Data Frames *groupby* method

Using "group by" method we can:

- Split the data into groups based on some criteria
- Calculate statistics (or apply a function) to each group
- Similar to dplyr() function in R

```
In [ ]: #Group data using rank  
df_rank = df.groupby(['rank'])
```

```
In [ ]: #Calculate mean value for each numeric column per each group  
df_rank.mean()
```

	phd	service	salary
rank			
<b>AssocProf</b>	15.076923	11.307692	91786.230769
<b>AsstProf</b>	5.052632	2.210526	81362.789474
<b>Prof</b>	27.065217	21.413043	123624.804348

# Data Frames *groupby* method

Once groupby object is created we can calculate various statistics for each group:

```
In [ ]: #Calculate mean salary for each professor rank:  
df.groupby('rank')[['salary']].mean()
```

rank	salary
AssocProf	91786.230769
AsstProf	81362.789474
Prof	123624.804348

Note: If single brackets are used to specify the column (e.g. salary), then the output is Pandas Series object. When double brackets are used the output is a Data Frame

# Data Frames *groupby* method

*groupby* performance notes:

- no grouping/splitting occurs until it's needed. Creating the *groupby* object only verifies that you have passed a valid mapping
- by default the group keys are sorted during the *groupby* operation. You may want to pass `sort=False` for potential speedup:

```
In [ ]: #Calculate mean salary for each professor rank:  
df.groupby(['rank'], sort=False) [['salary']].mean()
```

filtering

# Data Frame: filtering

To subset the data we can apply Boolean indexing. This indexing is commonly known as a filter. For example if we want to subset the rows in which the salary value is greater than \$120K:

```
In [ ]: #Calculate mean salary for each professor rank:  
df_sub = df[ df['salary'] > 120000 ]
```

Any Boolean operator can be used to subset the data:

- > greater;    >= greater or equal;
- < less;       <= less or equal;
- == equal;      != not equal;

```
In [ ]: #Select only those rows that contain female professors:  
df_f = df[ df['sex'] == 'Female' ]
```

slicing

# Data Frames: Slicing

There are a number of ways to subset the Data Frame:

- one or more columns
- one or more rows
- a subset of rows and columns

Rows and columns can be selected by their position or label

# Data Frames: Slicing

When selecting one column, it is possible to use single set of brackets, but the resulting object will be a Series (not a DataFrame):

```
In [ ]: #Select column salary:  
        df['salary']
```

When we need to select more than one column and/or make the output to be a DataFrame, we should use double brackets:

```
In [ ]: #Select column salary:  
        df[['rank', 'salary']]
```

# Data Frames: Selecting rows

If we need to select a range of rows, we can specify the range using ":"

```
In [ ]: #Select rows by their position:  
df[10:20]
```

Notice that the first row has a position 0, and the last value in the range is omitted: So for 0:10 range the first 10 rows are returned with the positions starting with 0 and ending with 9

# Data Frames: method loc

If we need to select a range of rows, using their labels we can use method loc:

```
In [ ]: #Select rows by their labels:  
df_sub.loc[10:20, ['rank', 'sex', 'salary']]
```

```
Out[ ]:
```

	rank	sex	salary
10	Prof	Male	128250
11	Prof	Male	134778
13	Prof	Male	162200
14	Prof	Male	153750
15	Prof	Male	150480
19	Prof	Male	150500

# Data Frames: method iloc

If we need to select a range of rows and/or columns, using their positions we can use method iloc:

```
In [ ]: #Select rows by their labels:  
df_sub.iloc[10:20, [0, 3, 4, 5]]
```

Out[ ]:

	rank	service	sex	salary
26	Prof	19	Male	148750
27	Prof	43	Male	155865
29	Prof	20	Male	123683
31	Prof	21	Male	155750
35	Prof	23	Male	126933
36	Prof	45	Male	146856
39	Prof	18	Female	129000
40	Prof	36	Female	137000
44	Prof	19	Female	151768
45	Prof	25	Female	140096

# Data Frames: method iloc (summary)

```
df.iloc[0]    # First row of a data frame  
df.iloc[i]    #(i+1)th row  
df.iloc[-1]   # Last row
```

```
df.iloc[:, 0]  # First column  
df.iloc[:, -1] # Last column
```

```
df.iloc[0:7]      #First 7 rows  
df.iloc[:, 0:2]    #First 2 columns  
df.iloc[1:3, 0:2]  #Second through third rows and first 2 columns  
df.iloc[[0,5], [1,3]] #1st and 6th rows and 2nd and 4th columns
```

sorting

# Data Frames: Sorting

We can sort the data by a value in the column. By default the sorting will occur in ascending order and a new data frame is return.

```
In [ ]: # Create a new data frame from the original sorted by the column Salary  
df_sorted = df.sort_values( by ='service')  
df_sorted.head()
```

```
Out[ ]:
```

	rank	discipline	phd	service	sex	salary
55	AsstProf	A	2	0	Female	72500
23	AsstProf	A	2	0	Male	85000
43	AsstProf	B	5	0	Female	77000
17	AsstProf	B	4	0	Male	92000
12	AsstProf	B	1	0	Male	88000

# Data Frames: Sorting

We can sort the data using 2 or more columns:

```
In [ ]: df_sorted = df.sort_values( by =['service', 'salary'], ascending = [True, False])
df_sorted.head(10)
```

Out[ ]:

	rank	discipline	phd	service	sex	salary
52	Prof	A	12	0	Female	105000
17	AsstProf	B	4	0	Male	92000
12	AsstProf	B	1	0	Male	88000
23	AsstProf	A	2	0	Male	85000
43	AsstProf	B	5	0	Female	77000
55	AsstProf	A	2	0	Female	72500
57	AsstProf	A	3	1	Female	72500
28	AsstProf	B	7	2	Male	91300
42	AsstProf	B	4	2	Female	80225
68	AsstProf	A	4	2	Female	77500

dropna

# Missing Values

Missing values are marked as NaN

```
In [ ]: # Read a dataset with missing values
flights = pd.read_csv("http://rcs.bu.edu/examples/python/data_analysis/flights.csv")
```

```
In [ ]: # Select the rows that have at least one missing value
flights[flights.isnull().any(axis=1)].head()
```

```
Out[ ]:
```

	year	month	day	dep_time	dep_delay	arr_time	arr_delay	carrier	tailnum	flight	origin	dest	air_time	distance	hour	minute
330	2013	1	1	1807.0	29.0	2251.0	NaN	UA	N31412	1228	EWR	SAN	NaN	2425	18.0	7.0
403	2013	1	1	NaN	NaN	NaN	NaN	AA	N3EHA	791	LGA	DFW	NaN	1389	NaN	NaN
404	2013	1	1	NaN	NaN	NaN	NaN	AA	N3EVAA	1925	LGA	MIA	NaN	1096	NaN	NaN
855	2013	1	2	2145.0	16.0	NaN	NaN	UA	N12221	1299	EWR	RSW	NaN	1068	21.0	45.0
858	2013	1	2	NaN	NaN	NaN	NaN	AA	NaN	133	JFK	LAX	NaN	2475	NaN	NaN

# Missing Values

There are a number of methods to deal with missing values in the data frame:

df.method()	description
dropna()	Drop missing observations
dropna(how='all')	Drop observations where all cells is NA
dropna(axis=1, how='all')	Drop column if all the values are missing
dropna(thresh = 5)	Drop rows that contain less than 5 non-missing values
fillna(0)	Replace missing values with zeros
isnull()	returns True if the value is missing
notnull()	Returns True for non-missing values

# pandas.DataFrame.dropna

DataFrame.`dropna(axis=0, how='any', thresh=None, subset=None, inplace=False)`

[source]

Return object with labels on given axis omitted where alternately any or all of the data are missing

## Parameters:

**axis** : {0 or ‘index’, 1 or ‘columns’}, or tuple/list thereof

Pass tuple or list to drop on multiple axes

**how** : {'any', 'all'}

- any : if any NA values are present, drop that label
- all : if all values are NA, drop that label

**thresh** : int, default None

int value : require that many non-NA values

**subset** : array-like

Labels along other axis to consider, e.g. if you are dropping rows these would be a list of columns to include

**inplace** : boolean, default False

If True, do operation inplace and return None.

## Returns:

**dropped** : DataFrame

```
7 import pandas as pd
8 import numpy as np
9 df = pd.DataFrame([[np.nan, 2, np.nan, 0], [3, 4, np.nan, 1],
10                      [np.nan, np.nan, np.nan, 5]],
11                      columns=list('ABCD'))
12 ...
13      A    B    C    D
14 0   NaN  2.0  NaN  0
15 1   3.0  4.0  NaN  1
16 2   NaN  NaN  NaN  5
17 ...
18 df.dropna(axis=1, how='all')
19 #drop the column only if all column values are NaN
20 ...
21      A    B    D
22 0   NaN  2.0  0
23 1   3.0  4.0  1
24 2   NaN  NaN  5
25 ...
26 df.dropna(axis=1, how='any')
27 #drop the column if at least one NaN cell, so only the last column remains
28 df.dropna(thresh=2)
29 #Keep only the rows with at least 2 non-na values
30 ...
31      A    B    C    D
32 0   NaN  2.0  NaN  0
33 1   3.0  4.0  NaN  1
34 ...
```

# pandas.Series.dropna

`Series.dropna(axis=0, inplace=False, **kwargs)`

Return Series without null values

**Returns:**

**valid** : Series

**inplace** : boolean, default *False*

Do operation in place.

```
from pandas import Series
import numpy as np

s1 = Series([1, 2, np.NaN])
s = s1.dropna(axis=0, inplace=False)
...
0    1
1    2
...
#Remove null values from series
```

# Missing Values

- When summing the data, missing values will be treated as zero
- If all values are missing, the sum will be equal to NaN
- `cumsum()` and `cumprod()` methods ignore missing values but preserve them in the resulting arrays
- Missing values in GroupBy method are excluded (just like in R)
- Many descriptive statistics methods have `skipna` option to control if missing data should be excluded . This value is set to `True` by default (unlike R)

aggregation

# Aggregation Functions in Pandas

Aggregation - computing a summary statistic about each group, i.e.

- compute group sums or means
- compute group sizes/counts

Common aggregation functions:

min, max

count, sum, prod

mean, median, mode, mad

std, var

# Aggregation Functions in Pandas

`agg()` method are useful when multiple statistics are computed per column:

```
In [ ]: flights[['dep_delay', 'arr_delay']].agg(['min', 'mean', 'max'])
```

Out[ ]:

	dep_delay	arr_delay
<b>min</b>	-16.000000	-62.000000
<b>mean</b>	9.384302	2.298675
<b>max</b>	351.000000	389.000000

## pandas.core.groupby.DataFrameGroupBy.aggregate

GroupBy.aggregate(func, \*args, \*\*kwargs)

[source]

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Sat Dec 30 15:10:21 2017
4 @author: Liu Jiali
5 """
6 import pandas as pd
7 df = pd.DataFrame({'A': [1,1,1,2,3,3,4,5], 'B': [1,1,1,3,4,4,5,5], 'C': [1,1,1,1,1,2,3,4], 'D': [5,6,7,8,9,1,1,1]})  
8 '''data looks like:  
9   A   B   C   D  
10  0   1   1   5  
11  1   1   1   6  
12  2   1   1   7  
13  3   2   3   8  
14  4   3   4   1  
15  5   3   4   2  
16  6   4   5   3  
17  7   5   5   4  
18  ...
19
20 df.groupby(['A','B']).aggregate(['min','max'])
21 # output:
22 #      C          D
23 #      min max min max
24 # A B
25 # 1 1   1   1   5   7
26 # 2 3   1   1   8   8
27 # 3 4   1   2   1   9
28 # 4 5   3   3   1   1
29 # 5 5   4   4   1   1
```

# Basic Descriptive Statistics

df.method()	description
describe	Basic statistics (count, mean, std, min, quantiles, max)
min, max	Minimum and maximum values
mean, median, mode	Arithmetic average, median and mode
var, std	Variance and standard deviation
sem	Standard error of mean
skew	Sample skewness
kurt	kurtosis

pd.DataFrame.fillna

# pandas.DataFrame.fillna

`DataFrame.fillna(value=None, method=None, axis=None, inplace=False, limit=None, downcast=None, **kwargs)`

Fill NA/NaN values using the specified method

[\[source\]](#)

## Parameters:

**value** : scalar, dict, Series, or DataFrame

Value to use to fill holes (e.g. 0), alternately a dict/Series/DataFrame of values specifying which value to use for each index (for a Series) or column (for a DataFrame). (values not in the dict/Series/DataFrame will not be filled). This value cannot be a list.

**method** : {'backfill', 'bfill', 'pad', 'ffill', None}, default None

Method to use for filling holes in reindexed Series pad / ffill: propagate last valid observation forward to next valid backfill / bfill: use NEXT valid observation to fill gap

**axis** : {0 or 'index', 1 or 'columns'}

**inplace** : boolean, default False

If True, fill in place. Note: this will modify any other views on this object, (e.g. a no-copy slice for a column in a DataFrame).

**limit** : int, default None

If method is specified, this is the maximum number of consecutive NaN values to forward/backward fill. In other words, if there is a gap with more than this number of consecutive NaNs, it will only be partially filled. If method is not specified, this is the maximum number of entries along the entire axis where NaNs will be filled. Must be greater than 0 if not None.

**downcast** : dict, default is None

a dict of item->dtype of what to downcast if possible, or the string 'infer' which will try to downcast to an appropriate equal type (e.g. float64 to int64 if possible)

```
7 import pandas as pd
8 import numpy as np
9 df = pd.DataFrame([[np.nan, 2, np.nan, 0],
10                     [3, 4, np.nan, 1],
11                     [np.nan, np.nan, np.nan, 5],
12                     [np.nan, 3, np.nan, 4]],
13                     columns=list('ABCD'))
14 df.fillna(0)
15 #Replace all NaN elements with 0s.
16 ...
17      A    B    C    D
18 0    0.0  2.0  0.0  0
19 1    3.0  4.0  0.0  1
20 2    0.0  0.0  0.0  5
21 3    0.0  3.0  0.0  4
22 ...
23 df.fillna(method='ffill')
24 #propagate non-null values forward
25 ...
26      A    B    C    D
27 0    NaN  2.0  NaN  0
28 1    3.0  4.0  NaN  1
29 2    3.0  4.0  NaN  5
30 3    3.0  3.0  NaN  4
31 ...
```

```
32 values = {'A': 0, 'B': 1, 'C': 2, 'D': 3}
33 df.fillna(value=values)
34 #Replace all NaN elements in column 'A', 'B', 'C', and 'D', with 0, 1, 2, and
35 #3 respectively.
36 ...
37      A    B    C    D
38 0    0.0  2.0  2.0  0
39 1    3.0  4.0  2.0  1
40 2    0.0  1.0  2.0  5
41 3    0.0  3.0  2.0  4
42 ...
43 df.fillna(value=values, limit=1)
44 #Only replace the first NaN element.
45 ...
46      A    B    C    D
47 0    0.0  2.0  2.0  0
48 1    3.0  4.0  NaN  1
49 2    NaN  1.0  NaN  5
50 3    NaN  3.0  NaN  4
51 ...
```

```
pd.DataFrame.drop_duplicates
```

# pandas.DataFrame.drop\_duplicates

`DataFrame.drop_duplicates(subset=None, keep='first', inplace=False)`

[\[source\]](#)

Return DataFrame with duplicate rows removed, optionally only considering certain columns

## Parameters:

**subset** : *column label or sequence of labels, optional*

Only consider certain columns for identifying duplicates, by default use all of the columns

**keep** : {‘first’, ‘last’, *False*}, default ‘first’

- *first* : Drop duplicates except for the first occurrence.
- *last* : Drop duplicates except for the last occurrence.
- *False* : Drop all duplicates.

**inplace** : *boolean, default False*

Whether to drop duplicates in place or to return a copy

## Returns:

**deduplicated** : *DataFrame*

```
# -*- coding: utf-8 -*-
"""
Created on Set Dec 23 2017
@author: Luo Hankun (3035084914)
"""

import pandas as pd

df = pd.DataFrame({'A': [1,1,2,3,5,6,7,8], 'B': [1,1,6,7,8,9,10,11]})

df.drop_duplicates()
"""

Output:
   A   B
0  1   1
2  2   6
3  3   7
4  5   8
5  6   9
6  7  10
7  8  11

Removed row 1 (duplicate)
"""
```

pd.DataFrame.T

# pandas.DataFrame.T

DataFrame.T

Transpose index and columns

```
7 import pandas as pd
8 df = pd.DataFrame({'A': {0: 'a', 1: 'b', 2: 'c'},
9                     'B': {0: 1, 1: 3, 2: 5},
10                    'C': {0: 2, 1: 4, 2: 6}})
11 ...
12    A   B   C
13 0  a  1  2
14 1  b  3  4
15 2  c  5  6
16 ...
17 df.T
18 #Transpose index and columns
19 ...
20    0   1   2
21 A  a  b  c
22 B  1  3  5
23 C  2  4  6
24 ...
```

pd.unique

# pandas.unique

`pandas.unique(values)`

[source]

Hash table-based `unique`. Uniques are returned in order of appearance. This does NOT sort.

Significantly faster than `numpy.unique`. Includes NA values.

**Parameters:**

**values** : *1d array-like*

**Returns:**

*unique values*.

- If the input is an `Index`, the return is an `Index`
- If the input is a `Categorical dtype`, the return is a `Categorical`
- If the input is a `Series/ndarray`, the return will be an `ndarray`

**See also:** `pandas.Index.unique`, `pandas.Series.unique`

```
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3 """
4 Created on Thu Dec 21 17:26:38 2017
5
6 @author: DemiWu
7 """
8
9 import pandas as pd
10 df = pd.DataFrame({'A': [1, 3, 1, 4, 3], 'B': [3, 1, 4, 2, 2]})
11
12 #find unique values in column A
13 pd.unique(df.A)
14 '''the result is
15 array([1, 3, 4])
16 '''
```

pd.to\_datetime

# pandas.to\_datetime

```
pandas.to_datetime(arg, errors='raise', dayfirst=False, yearfirst=False, utc=None, box=True, format=None,  
exact=True, unit=None, infer_datetime_format=False, origin='unix', cache=False) [source]
```

Convert argument to datetime.

**arg** : integer, float, string, datetime, list, tuple, 1-d array, Series

New in version 0.18.1: or DataFrame/dict-like

**errors** : {'ignore', 'raise', 'coerce'}, default 'raise'

- If 'raise', then invalid parsing will raise an exception
- If 'coerce', then invalid parsing will be set as NaT
- If 'ignore', then invalid parsing will return the input

**dayfirst** : boolean, default False

Specify a date parse order if arg is str or its list-likes. If True, parses dates with the day first, eg 10/11/12 is parsed as 2012-11-10. Warning: dayfirst=True is not strict, but will prefer to parse with day first (this is a known bug, based on dateutil behavior).

**yearfirst** : boolean, default False

Specify a date parse order if arg is str or its list-likes.

- If True parses dates with the year first, eg 10/11/12 is parsed as 2010-11-12.
- If both dayfirst and yearfirst are True, yearfirst is preceded (same as dateutil).

Warning: yearfirst=True is not strict, but will prefer to parse with year first (this is a known bug, based on dateutil behavior).

New in version 0.16.1.

**utc** : boolean, default None

Return UTC DatetimeIndex if True (converting any tz-aware datetime.datetime objects as well).

**box** : boolean, default True

- If True returns a DatetimeIndex
- If False returns ndarray of values.

**format** : string, default None

strftime to parse time, eg "%d/%m/%Y", note that "%f" will parse all the way up to nanoseconds.

## Parameters:

**exact** : boolean, *True by default*

- If True, require an exact format match.
- If False, allow the format to match anywhere in the target string.

**unit** : string, *default 'ns'*

unit of the arg (D,s,ms,us,ns) denote the unit, which is an integer or float number. This will be based off the origin. Example, with unit='ms' and origin='unix' (the default), this would calculate the number of milliseconds to the unix epoch start.

**infer\_datetime\_format** : boolean, *default False*

If True and no *format* is given, attempt to infer the format of the datetime strings, and if it can be inferred, switch to a faster method of parsing them. In some cases this can increase the parsing speed by ~5-10x.

**origin** : scalar, *default is 'unix'*

Define the reference date. The numeric values would be parsed as number of units (defined by *unit*) since this reference date.

- If 'unix' (or POSIX) time; origin is set to 1970-01-01.
- If 'julian', unit must be 'D', and origin is set to beginning of Julian Calendar. Julian day number 0 is assigned to the day starting at noon on January 1, 4713 BC.
- If Timestamp convertible, origin is set to Timestamp identified by origin.

*New in version 0.20.0.*

**cache** : boolean, *default False*

If True, use a cache of unique, converted dates to apply the datetime conversion. May produce significant speed-up when parsing duplicate date strings, especially ones with timezone offsets.

*New in version 0.23.0.*

**ret** : datetime if parsing succeeded.

Return type depends on input:

- list-like: DatetimeIndex
- Series: Series of datetime64 dtype
- scalar: Timestamp

In case when it is not possible to return designated types (e.g. when any element of input is before Timestamp.min or after Timestamp.max) return will have datetime.datetime type (or corresponding array/Series).

**Returns:**

```
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3 """
4 Created on Fri Dec 22 16:53:43 2017
5
6 @author: DemiWu
7 """
8
9 import pandas as pd
10 dt = ['2017-01-05 2:30:00 PM', 'Jan 5, 2017 14:30:00', '01/05/2016',
11       '2017.01.05', '2017/01/05', '20170105']
12
13 #convert dt to datetime format
14 pd.to_datetime(dt)
15 """the result is
16 DatetimeIndex(['2017-01-05 14:30:00', '2017-01-05 14:30:00',
17                 '2016-01-05 00:00:00', '2017-01-05 00:00:00',
18                 '2017-01-05 00:00:00', '2017-01-05 00:00:00'],
19                 dtype='datetime64[ns]', freq=None)
20 """
```

```
31 #about "yearfirst" argument
32 pd.to_datetime('10/11/12')
33 '''the result is
34 Timestamp('2012-10-11 00:00:00')
35 '''
36 pd.to_datetime('10/11/12', yearfirst=True)
37 '''the result is
38 Timestamp('2010-11-12 00:00:00')
39 '''
40 #note that when yearfirst is set to true, the first number in the input is recognized as year
41
42 #custom datetime format of the input
43 pd.to_datetime('2017$01$05')
44 '''ValueError: Unknown string format'''
45 pd.to_datetime('2017$01$05', format='%Y$%m$%d')
46 '''the result is
47 Timestamp('2017-01-05 00:00:00')
48 '''
49 pd.to_datetime('2017/01/05', infer_datetime_format=True)
50 '''the result is Timestamp('2017-01-05 00:00:00')'''
51
52 #handle invalid dates
53 pd.to_datetime(['2017-01-05', 'Jan 6, 2017', 'abc'])
54 '''ValueError: Unknown string format'''
55 pd.to_datetime(['2017-01-05', 'Jan 6, 2017', 'abc'], errors='ignore')
56 '''the result is
57 array(['2017-01-05', 'Jan 6, 2017', 'abc'], dtype=object)
58 '''
59 #note that the input is returned as original
60 pd.to_datetime(['2017-01-05', 'Jan 6, 2017', 'abc'], errors='coerce')
61 '''the result is
62 DatetimeIndex(['2017-01-05', '2017-01-06', 'NaT'], dtype='datetime64[ns]', freq=None)
63 '''
64 #note that 'abc' is returned as 'NaT'
65'''
```

```
66 #Unix time: the time has passed since the origin
67 pd.to_datetime(1501324478, unit='s')
68 '''the result is Timestamp('2017-07-29 10:34:38')'''
69 pd.to_datetime(1501324478000, unit='ms')
70 '''the result is Timestamp('2017-07-29 10:34:38')'''
71 #change the origin of unix time
72 pd.to_datetime(2458072.5, unit='D', origin='julian')
73 '''the result is Timestamp('2017-11-15 00:00:00')'''
74
75 #use utc time
76 pd.to_datetime(dt, utc=True)
77 '''the result is
78 DatetimeIndex(['2017-01-05 14:30:00+00:00', '2017-01-05 14:30:00+00:00',
79                 '2016-01-05 00:00:00+00:00', '2017-01-05 00:00:00+00:00',
80                 '2017-01-05 00:00:00+00:00', '2017-01-05 00:00:00+00:00'],
81                 dtype='datetime64[ns, UTC]', freq=None)
82 '''
83
84 #change types of return
85 pd.to_datetime(dt, box=False)
86 '''the result is
87 array(['2017-01-05T14:30:00.000000000', '2017-01-05T14:30:00.000000000',
88        '2016-01-05T00:00:00.000000000', '2017-01-05T00:00:00.000000000',
89        '2017-01-05T00:00:00.000000000', '2017-01-05T00:00:00.000000000'],
90        dtype='datetime64[ns]')
91 #note that the return is array instead of datetimeindex
```

pd.DataFrame.merge

# pandas.DataFrame.merge

```
DataFrame.merge(right, how='inner', on=None, left_on=None, right_on=None, left_index=False, right_index=False,  
sort=False, suffixes=('_x', '_y'), copy=True, indicator=False, validate=None) [source]
```

Merge DataFrame objects by performing a database-style join operation by columns or indexes.

If joining columns on columns, the DataFrame indexes *will be ignored*. Otherwise if joining indexes on indexes or indexes on a column or columns, the index will be passed on.

**right** : *DataFrame*

**how** : {'left', 'right', 'outer', 'inner'}, default 'inner'

- left: use only keys from left frame, similar to a SQL left outer join; preserve key order
- right: use only keys from right frame, similar to a SQL right outer join; preserve key order
- outer: use union of keys from both frames, similar to a SQL full outer join; sort keys lexicographically
- inner: use intersection of keys from both frames, similar to a SQL inner join; preserve the order of the left keys

**on** : *label or list*

Column or index level names to join on. These must be found in both DataFrames. If *on* is None and not merging on indexes then this defaults to the intersection of the columns in both DataFrames.

**left\_on** : *label or list, or array-like*

Column or index level names to join on in the left DataFrame. Can also be an array or list of arrays of the length of the left DataFrame. These arrays are treated as if they are columns.

**right\_on** : *label or list, or array-like*

Column or index level names to join on in the right DataFrame. Can also be an array or list of arrays of the length of the right DataFrame. These arrays are treated as if they are columns.

**Parameters:**

**left\_index** : boolean, default False  
Use the index from the left DataFrame as the join key(s). If it is a MultiIndex, the number of keys in the other DataFrame (either the index or a number of columns) must match the number of levels

**right\_index** : boolean, default False  
Use the index from the right DataFrame as the join key. Same caveats as `left_index`

**sort** : boolean, default False  
Sort the join keys lexicographically in the result DataFrame. If False, the order of the join keys depends on the join type (how keyword)

**suffixes** : 2-length sequence (tuple, list, ...)  
Suffix to apply to overlapping column names in the left and right side, respectively

**copy** : boolean, default True  
If False, do not copy data unnecessarily

**indicator** : boolean or string, default False  
If True, adds a column to output DataFrame called “\_merge” with information on the source of each row. If string, column with information on source of each row will be added to output DataFrame, and column will be named value of string. Information column is Categorical-type and takes on a value of “left\_only” for observations whose merge key only appears in ‘left’ DataFrame, “right\_only” for observations whose merge key only appears in ‘right’ DataFrame, and “both” if the observation’s merge key is found in both.

**validate** : string, default None  
If specified, checks if merge is of specified type.

- “one\_to\_one” or “1:1”: check if merge keys are unique in both left and right datasets.
- “one\_to\_many” or “1:m”: check if merge keys are unique in left dataset.
- “many\_to\_one” or “m:1”: check if merge keys are unique in right dataset.
- “many\_to\_many” or “m:m”: allowed, but does not result in checks.

*New in version 0.21.0.*

**Returns:****merged** : DataFrame

The output type will be same as ‘left’, if it is a subclass of DataFrame.

```
8 import pandas as pd
9
10 caller = pd.DataFrame({'key': ['K0', 'K1', 'K2', 'K3', 'K4', 'K5'],
11                         'A': ['A0', 'A1', 'A2', 'A3', 'A4', 'A5']})
12 ...
13     A key
14 0   A0   K0
15 1   A1   K1
16 2   A2   K2
17 3   A3   K3
18 4   A4   K4
19 5   A5   K5
20 ...
21 other = pd.DataFrame({'key': ['K0', 'K1', 'K2'],
22                       'B': ['B0', 'B1', 'B2']})
23 ...
24     B key
25 0   B0   K0
26 1   B1   K1
27 2   B2   K2
28 ...
29 caller.merge(other, left_on='key', right_on='key', how='inner')
30 #left_on and right_on: Field names to join
31 #inner only keeps the intersection by 'key'
32 ...
33     A key    B
34 0   A0   K0   B0
35 1   A1   K1   B1
36 2   A2   K2   B2
37 ...
```

```
38 """
39 how : {'left', 'right', 'outer', 'inner'}, default 'inner'
40 left: use only keys from left frame, similar to a SQL left outer join; preserve
41     key order
42 right: use only keys from right frame, similar to a SQL right outer join;
43     preserve key order
44 outer: use union of keys from both frames, similar to a SQL full outer join;
45     sort keys lexicographically
46 inner: use intersection of keys from both frames, similar to a SQL inner join;
47     preserve the order of the left keys
48 """
49 caller.merge(other, on='key', how='outer',sort=1)
50 #on can be used when both df have the same key column
51 #outer join keeps the union by 'key'
52 #sort=1: Sort the join keys Lexicographically in the result DataFrame. If
53 #False, the order of the join keys depends on the join type
54 """
55      A   key      B
56 0    A0    K0    B0
57 1    A1    K1    B1
58 2    A2    K2    B2
59 3    A3    K3    NaN
60 4    A4    K4    NaN
61 5    A5    K5    NaN
62 """
63 caller.merge(other, on='key', how='right',sort=1)
64 #use only keys from right frame, similar to a SQL right outer join; preserve
65 #key order
66 """
67      A   key      B
68 0    A0    K0    B0
69 1    A1    K1    B1
70 2    A2    K2    B2
71 """
```

pd.concat

```
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3 """
4 Created on Mon Dec 25 15:06:38 2017
5
6 @author: jessicaleung
7 """
8
9 import pandas as pd
10
11 s1 = pd.Series(['a', 'b'])
12 s2 = pd.Series(['c', 'd'])
13
14 pd.concat([s1, s2])
15 # Output:
16 # 0    a
17 # 1    b
18 # 0    c
19 # 1    d
20 # dtype: object
21
22 pd.concat([s1, s2], axis=1)
23 # Output:
24 #   0  1
25 # 0  a  c
26 # 1  b  d
27
28 """ Set axis=1 to concatenate along columns. """
29
30 pd.concat([s1, s2], ignore_index=True)
31 # Output:
32 # 0    a
33 # 1    b
34 # 2    c
35 # 3    d
36 # dtype: object
37
38 """ Set ignore_index=True to assign new index values. It is by default False to
39 use the index values along the concatenation axis. """
40
```

```
40
41 pd.concat([s1, s2], keys=['s1', 's2'])
42 # Output:
43 # s1  0    a
44 #      1    b
45 # s2  0    c
46 #      1    d
47 # dtype: object
48
49 ''' Set keys to construct hierarchical index as the outermost level. '''
50
51 pd.concat([s1, s2], keys=['s1', 's2'], names=['Series name', 'Row ID'])
52 # Output:
53 # Series name  Row ID
54 # s1           0        a
55 #                 1        b
56 # s2           0        c
57 #                 1        d
58 # dtype: object
59
60 ''' Add Names for the levels in the resulting hierarchical index. '''
```

```
61
62 s3 = pd.Series([1,2,3], index=['a','b','c'])
63 ...
64 s3 is like:
65 a    1
66 b    2
67 c    3
68 dtype: int64
69 ...
70
71 s4 = pd.Series([5,6,7], index=['a','c','d'])
72 ...
73 s4 is like:
74 a    5
75 c    6
76 d    7
77 dtype: int64
78 ...
79
80 pd.concat([s3,s4], axis=1)
81 # Output:
82 #      0    1
83 # a  1.0  5.0
84 # b  2.0  NaN
85 # c  3.0  6.0
86 # d  NaN  7.0
87
88 pd.concat([s3,s4], axis=1, join_axes=[['a','b','c']])
89 # Output:
90 #      0    1
91 # a  1  5.0
92 # b  2  NaN
93 # c  3  6.0
94
95 """ Specify indexes to use for the other n - 1 axes instead of performing inner/outer set logic. """
96
```

```
96
97 df1 = pd.DataFrame([['a', 1], ['b', 2]], columns=['letter', 'number'])
98 """
99 df1 is like:
100    letter  number
101 0      a      1
102 1      b      2
103 """
104
105 df2 = pd.DataFrame([['c', 3], ['d', 4]], columns=['letter', 'number'])
106 """
107 df2 is like:
108    letter  number
109 0      c      3
110 1      d      4
111 """
112
113 pd.concat([df1, df2])
114 # Output:
115 #    letter  number
116 # 0      a      1
117 # 1      b      2
118 # 0      c      3
119 # 1      d      4
120
121 pd.concat([df1, df2], ignore_index=True)
122 # Output:
123 #    letter  number
124 # 0      a      1
125 # 1      b      2
126 # 2      c      3
127 # 3      d      4
128
129 new_df=pd.concat([df1, df2], axis=1, keys=['lv1','lv2'], names=['upper', 'lower'], levels=[['lv1','lv2','lv3']])
130 # Output:
131 # upper    lv1          lv2
132 # lower letter number letter number
133 # 0        a      1      c      3
134 # 1        b      2      d      4
135
136 new_df.columns
137 # Output:
138 # MultiIndex(levels=[['lv1', 'lv2', 'lv3'], ['letter', 'number']],
139 #             labels=[[0, 0, 1, 1], [0, 1, 0, 1]],
140 #             names=['upper', 'lower'])
141
142 """ adding levels to the parameter specifies levels to construct MultiIndex. """
143
```

```
143
144 df3 = pd.DataFrame([['c', 3, 'cat'], ['d', 4, 'dog']], columns=['letter', 'number', 'animal'])
145 """
146 df3 is like:
147   letter  number animal
148 0      c        3    cat
149 1      d        4    dog
150 """
151
152 pd.concat([df1, df2, df3])
153 # Output:
154 #   animal letter  number
155 # 0      NaN      a      1
156 # 1      NaN      b      2
157 # 0      NaN      c      3
158 # 1      NaN      d      4
159 # 0      cat      c      3
160 # 1      dog      d      4
161
162 pd.concat([df1, df2, df3], join='inner')
163 # Output:
164 #   letter  number
165 # 0      a      1
166 # 1      b      2
167 # 0      c      3
168 # 1      d      4
169 # 0      c      3
170 # 1      d      4
171
172 """ Column with NaN entries are excluded when join='inner' """

```

```
173
174 df5 = pd.DataFrame([1,5], index=['a','b'])
175 """
176 df5 is like:
177     0
178 a   1
179 b   5
180 """
181
182 df6 = pd.DataFrame([2,6], index=['a','c'])
183 """
184 df6 is like:
185     0
186 a   2
187 c   6
188 """
189
190 pd.concat([df5, df6])
191 # Output:
192 #     0
193 # a   1
194 # b   5
195 # a   2
196 # c   6
197
198 pd.concat([df5, df6], verify_integrity=True)
199 # Output:
200 # ValueError: Indexes have overlapping values: ['a']
201
202 """ Check whether the new concatenated axis contains duplicates. """
203
```

pd.merge

# Merge, join, and concatenate

- <https://pandas.pydata.org/pandas-docs/stable/merging.html>

```
pd.merge(left, right, how='inner', on=None, left_on=None, right_on=None,
         left_index=False, right_index=False, sort=True,
         suffixes=('_x', '_y'), copy=True, indicator=False,
         validate=None)
```

- `left`: A DataFrame object
- `right`: Another DataFrame object
- `on`: Columns (names) to join on. Must be found in both the left and right DataFrame objects. If not passed and `left_index` and `right_index` are `False`, the intersection of the columns in the DataFrames will be inferred to be the join keys
- `left_on`: Columns from the left DataFrame to use as keys. Can either be column names or arrays with length equal to the length of the DataFrame
- `right_on`: Columns from the right DataFrame to use as keys. Can either be column names or arrays with length equal to the length of the DataFrame
- `left_index`: If `True`, use the index (row labels) from the left DataFrame as its join key(s). In the case of a DataFrame with a MultiIndex (hierarchical), the number of levels must match the number of join keys from the right DataFrame
- `right_index`: Same usage as `left_index` for the right DataFrame
- `how`: One of '`left`', '`right`', '`outer`', '`inner`'. Defaults to `inner`. See below for more detailed description of each method

- `sort`: Sort the result DataFrame by the join keys in lexicographical order. Defaults to `True`, setting to `False` will improve performance substantially in many cases
- `suffixes`: A tuple of string suffixes to apply to overlapping columns. Defaults to `(''_x'', ''_y'')`.
- `copy`: Always copy data (default `True`) from the passed DataFrame objects, even when reindexing is not necessary. Cannot be avoided in many cases but may improve performance / memory usage. The cases where copying can be avoided are somewhat pathological but this option is provided nonetheless.
- `indicator`: Add a column to the output DataFrame called `_merge` with information on the source of each row. `_merge` is Categorical-type and takes on a value of `left_only` for observations whose merge key only appears in 'left' DataFrame, `right_only` for observations whose merge key only appears in 'right' DataFrame, and `both` if the observation's merge key is found in both.

*New in version 0.17.0.*

- `validate` : string, default `None`. If specified, checks if merge is of specified type.
  - "one\_to\_one" or "1:1": checks if merge keys are unique in both left and right datasets.
  - "one\_to\_many" or "1:m": checks if merge keys are unique in left dataset.
  - "many\_to\_one" or "m:1": checks if merge keys are unique in right dataset.
  - "many\_to\_many" or "m:m": allowed, but does not result in checks.

```
pd.merge(left, right, how='inner', on=None, left_on=None, right_on=None,  
        left_index=False, right_index=False, sort=True,  
        suffixes=('_x', '_y'), copy=True, indicator=False,  
        validate=None)
```

**how : {'left', 'right', 'outer', 'inner'}, default 'inner'**

- **left**: use only keys from left frame, similar to a SQL left outer join; preserve key order
- **right**: use only keys from right frame, similar to a SQL right outer join; preserve key order
- **outer**: use union of keys from both frames, similar to a SQL full outer join; sort keys lexicographically
- **inner**: use intersection of keys from both frames, similar to a SQL inner join; preserve the order of the left keys

```
In [38]: left = pd.DataFrame({'key': ['K0', 'K1', 'K2', 'K3'],
....:                         'A': ['A0', 'A1', 'A2', 'A3'],
....:                         'B': ['B0', 'B1', 'B2', 'B3']})
....:

In [39]: right = pd.DataFrame({'key': ['K0', 'K1', 'K2', 'K3'],
....:                         'C': ['C0', 'C1', 'C2', 'C3'],
....:                         'D': ['D0', 'D1', 'D2', 'D3']})
....:

In [40]: result = pd.merge(left, right, on='key')
```

left			right			Result		
	A	B	key	C	D	key	A	B
0	A0	B0	K0	C0	D0	K0	A0	B0
1	A1	B1	K1	C1	D1	K1	A1	B1
2	A2	B2	K2	C2	D2	K2	A2	B2
3	A3	B3	K3	C3	D3	K3	A3	B3

```
In [46]: result = pd.merge(left, right, how='outer', on=['key1', 'key2'])
```

	A	B	key1	key2
0	A0	B0	K0	K0
1	A1	B1	K0	K1
2	A2	B2	K1	K0
3	A3	B3	K2	K1

	C	D	key1	key2
0	C0	D0	K0	K0
1	C1	D1	K1	K0
2	C2	D2	K1	K0
3	C3	D3	K2	K0

	A	B	key1	key2	C	D
0	A0	B0	K0	K0	C0	D0
1	A1	B1	K0	K1	NaN	NaN
2	A2	B2	K1	K0	C1	D1
3	A2	B2	K1	K0	C2	D2
4	A3	B3	K2	K1	NaN	NaN
5	NaN	NaN	K2	K0	C3	D3

```
....  
Created on Thu Nov 23 20:28:36 2017  
  
@author: wanglu  
....  
  
#import package  
import pandas as pd  
import os  
  
#define two DataFrame to be merged  
left = pd.DataFrame({'key' : ['a','b','c','d','e'],  
                     'value_1' : [1,2,3,4,5],  
                     'value_2' : ['I','II','III','IV','V']})  
  
right = pd.DataFrame({'key' : ['a','b','c','d'],  
                      'value_3' : ['A','B','C','D'],  
                      'value_4' : ['one','two','three','four']})  
  
#merge DataFrame 'right' to 'left' using variables in column 'key' in the left DataFrame as join keys  
# if how = 'outer', use union of keys from both frames  
result = pd.merge(left, right, how='outer', on='key')  
  
#save this DataFrame to a csv file  
data_path = '/Users/wanglu/Desktop/Pandas Functions.General Functions/Pandas Functions.General Functions.merge'  
result.to_csv(data_path + os.sep + 'General_Function_merge_result.csv')
```

---

	key	value_1	value_2	value_3	value_4
0	a	1	I	A	one
1	b	2	II	B	two
2	c	3	III	C	three
3	d	4	IV	D	four
4	e	5	V		

# pandas.merge

```
pandas.merge(left, right, how='inner', on=None, left_on=None, right_on=None, left_index=False, right_index=False,  
sort=False, suffixes=('_x', '_y'), copy=True, indicator=False, validate=None) [source]
```

Merge DataFrame objects by performing a database-style join operation by columns or indexes.

If joining columns on columns, the DataFrame indexes *will be ignored*. Otherwise if joining indexes on indexes or indexes on a column or columns, the index will be passed on.

**left** : *DataFrame*

**right** : *DataFrame*

**how** : {‘left’, ‘right’, ‘outer’, ‘inner’}, default ‘inner’

- left: use only keys from left frame, similar to a SQL left outer join; preserve key order
- right: use only keys from right frame, similar to a SQL right outer join; preserve key order
- outer: use union of keys from both frames, similar to a SQL full outer join; sort keys lexicographically
- inner: use intersection of keys from both frames, similar to a SQL inner join; preserve the order of the left keys

**on** : *label or list*

Field names to join on. Must be found in both DataFrames. If on is None and not merging on indexes, then it merges on the intersection of the columns by default.

**left\_on** : *label or list, or array-like*

Field names to join on in left DataFrame. Can be a vector or list of vectors of the length of the DataFrame to use a particular vector as the join key instead of columns

**right\_on** : *label or list, or array-like*

Field names to join on in right DataFrame or vector/list of vectors per left\_on docs

**Parameters:**

**left\_index** : boolean, default False

Use the index from the left DataFrame as the join key(s). If it is a MultiIndex, the number of keys in the other DataFrame (either the index or a number of columns) must match the number of levels

**right\_index** : boolean, default False

Use the index from the right DataFrame as the join key. Same caveats as `left_index`

**sort** : boolean, default False

Sort the join keys lexicographically in the result DataFrame. If False, the order of the join keys depends on the join type (how keyword)

**suffixes** : 2-length sequence (tuple, list, ...)

Suffix to apply to overlapping column names in the left and right side, respectively

**copy** : boolean, default True

If False, do not copy data unnecessarily

**indicator** : boolean or string, default False

If True, adds a column to output DataFrame called “\_merge” with information on the source of each row. If string, column with information on source of each row will be added to output DataFrame, and column will be named value of string. Information column is Categorical-type and takes on a value of “left\_only” for observations whose merge key only appears in ‘left’ DataFrame, “right\_only” for observations whose merge key only appears in ‘right’ DataFrame, and “both” if the observation’s merge key is found in both.

New in version 0.17.0.

**validate** : string, default None

If specified, checks if merge is of specified type.

- “one\_to\_one” or “1:1”: check if merge keys are unique in both left and right datasets.
- “one\_to\_many” or “1:m”: check if merge keys are unique in left dataset.
- “many\_to\_one” or “m:1”: check if merge keys are unique in right dataset.
- “many\_to\_many” or “m:m”: allowed, but does not result in checks.

New in version 0.21.0.

**Returns:**

**merged** : *DataFrame*

The output type will be the same as 'left', if it is a subclass of DataFrame.

**See also:** [merge\\_ordered](#), [merge\\_asof](#)

```
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3 """
4 Created on Wed Dec 20 16:44:28 2017
5
6 @author: DemiWu
7 """
8
9 import pandas as pd
10 df1 = pd.DataFrame({"city": ["new york", "chicago", "orlando", "baltimore"],
11                      "temperature": [21, 14, 35, 38]})
12 """df1 looks like
13      city  temperature
14 0    new york          21
15 1    chicago           14
16 2    orlando            35
17 3   baltimore           38
18 """
19
20 df2 = pd.DataFrame({"city": ["chicago", "new york", "san diego"],
21                      "humidity": [65, 68, 71]})
22 """data looks like
23      city  humidity
24 0    chicago          65
25 1    new york          68
26 2   san diego          71
27 """
```

```
29 #merge the two dataframes using only the cities that are in both dataframes
30 df3=pd.merge(df1,df2,on="city",how="inner")
31 '''the result is
32      city  temperature  humidity
33 0  new york          21        68
34 1  chicago           14        65
35 '''
36
37 #merge the two dataframes using the cities that are in either dataframes
38 df3=pd.merge(df1,df2,on="city",how="outer")
39 '''the result is
40      city  temperature  humidity
41 0  new york          21.0       68.0
42 1  chicago           14.0       65.0
43 2  orlando            35.0       NaN
44 3  baltimore          38.0       NaN
45 4  san diego          NaN        71.0
46 '''
47
48 #merge the two dataframes using the cities in df1
49 df3=pd.merge(df1,df2,on="city",how="left")
50 '''the result is
51      city  temperature  humidity
52 0  new york          21        68.0
53 1  chicago           14        65.0
54 2  orlando            35        NaN
55 3  baltimore          38        NaN
56 '''
57
58 #merge the two dataframes using the cities in df2
59 df3=pd.merge(df1,df2,on="city",how="right")
60 '''the result is
61      city  temperature  humidity
62 0  new york          21.0        68
63 1  chicago           14.0        65
64 2  san diego          NaN        71
65 '''
```

```
67 #add a column showing the sources of each row
68 df3=pd.merge(df1,df2,on="city",how="outer",indicator=True)
69 '''the result is
70      city  temperature  humidity      _merge
71 0  new york      21.0      68.0      both
72 1  chicago       14.0      65.0      both
73 2  orlando        35.0      NaN  left_only
74 3  baltimore      38.0      NaN  left_only
75 4  san diego      NaN      71.0 right_only
76 '''
77
78 #use the "left_on" and "right_on" argument to replace "on"
79 df3=pd.merge(df1,df2, left_on="city", right_on="city", how="outer",indicator=True)
80 '''the result is
81      city  temperature  humidity      _merge
82 0  new york      21.0      68.0      both
83 1  chicago       14.0      65.0      both
84 2  orlando        35.0      NaN  left_only
85 3  baltimore      38.0      NaN  left_only
86 4  san diego      NaN      71.0 right_only
87 '''
88 #the same as the previous result using "on" argument
89
```

```
89 #merge overlapping columns
90 df1 = pd.DataFrame({
91     "city": ["new york", "chicago", "orlando", "baltimore"],
92     "temperature": [21, 14, 35, 38],
93     "humidity": [65, 68, 71, 75]})
94 '''the data looks like
95      city    humidity    temperature
96 0  new york        65            21
97 1  chicago         68            14
98 2  orlando          71            35
99 3  baltimore        75            38
100 ...
101 ...
102 df2 = pd.DataFrame({
103     "city": ["chicago", "new york", "san diego"],
104     "temperature": [21, 14, 35],
105     "humidity": [65, 68, 71]})
106 '''the data looks like
107      city    humidity    temperature
108 0  chicago         65            21
109 1  new york        68            14
110 2  san diego       71            35
111 ...
112 df3= pd.merge(df1,df2, on="city", how="outer", suffixes=('_first','_second'))
113 '''the result is
114      city    humidity_first    temperature_first    humidity_second \
115 0  new york           65.0              21.0           68.0
116 1  chicago            68.0              14.0           65.0
117 2  orlando             71.0              35.0            NaN
118 3  baltimore           75.0              38.0            NaN
119 4  san diego           NaN                NaN           71.0
120 ...
121      temperature_second
122 0                  14.0
123 1                  21.0
124 2                  NaN
125 3                  NaN
126 4                  35.0
127 ...
```

```
129 #use the index as the join key
130 df1 = pd.DataFrame({
131     "city": [1,2,3],
132     "temperature": [21,14,35],})
133 df1.set_index('city',inplace=True)
134 '''the data looks like
135     temperature
136 city
137 1          21
138 2          14
139 3          35
140 ...
141 df2 = pd.DataFrame({
142     "city": [2,3,4],
143     "humidity": [65,68,75],})
144 df2.set_index('city',inplace=True)
145 '''the data looks like
146     humidity
147 city
148 2          65
149 3          68
150 4          75
151 ...
152 df3= pd.merge(df1,df2, left_index=True, right_index=True, how="outer")
153 '''the result is
154     temperature  humidity
155 city
156 1          21.0      NaN
157 2          14.0      65.0
158 3          35.0      68.0
159 4          NaN        75.0
160 ...
161
```

```
161
162 #about the sort argument
163 df1 = pd.DataFrame({
164     "city": ["new york", "chicago", "orlando"],
165     "temperature": [21, 14, 35], })
166 df1.set_index('city', inplace=True)
167 df2 = pd.DataFrame({
168     "city": ["chicago", "new york", "orlando"],
169     "humidity": [65, 68, 75], })
170 df2.set_index('city', inplace=True)
171 df3= pd.merge(df1,df2, left_index=True, right_index=True, how="outer", sort=True)
172 '''the result is
173             temperature    humidity
174 city
175 chicago            14          65
176 new york           21          68
177 orlando            35          75
178 '''
179 #the join keys are sorted lexicographically
180
```

pd.merge\_asof

# pandas.merge\_asof

```
pandas.merge_asof(left, right, on=None, left_on=None, right_on=None, left_index=False, right_index=False,  
by=None, left_by=None, right_by=None, suffixes=('_x', '_y'), tolerance=None, allow_exact_matches=True,  
direction='backward')
```

[\[source\]](#)

Perform an asof merge. This is similar to a left-join except that we match on nearest key rather than equal keys.

Both DataFrames must be sorted by the key.

For each row in the left DataFrame:

- A “backward” search selects the last row in the right DataFrame whose ‘on’ key is less than or equal to the left’s key.
- A “forward” search selects the first row in the right DataFrame whose ‘on’ key is greater than or equal to the left’s key.
- A “nearest” search selects the row in the right DataFrame whose ‘on’ key is closest in absolute distance to the left’s key.

The default is “backward” and is compatible in versions below 0.20.0. The direction parameter was added in version 0.20.0 and introduces “forward” and “nearest”.

Optionally match on equivalent keys with ‘by’ before searching with ‘on’.

*New in version 0.19.0.*

**Parameters:**

<b>left</b> : <i>DataFrame</i>	
<b>right</b> : <i>DataFrame</i>	
<b>on</b> : <i>label</i>	Field name to join on. Must be found in both DataFrames. The data MUST be ordered. Furthermore this must be a numeric column, such as datetimelike, integer, or float. On or left_on/right_on must be given.
<b>left_on</b> : <i>label</i>	Field name to join on in left DataFrame.
<b>right_on</b> : <i>label</i>	Field name to join on in right DataFrame.
<b>left_index</b> : <i>boolean</i>	Use the index of the left DataFrame as the join key. <i>New in version 0.19.2.</i>
<b>right_index</b> : <i>boolean</i>	Use the index of the right DataFrame as the join key. <i>New in version 0.19.2.</i>
<b>by</b> : <i>column name or list of column names</i>	Match on these columns before performing merge operation.
<b>left_by</b> : <i>column name</i>	Field names to match on in the left DataFrame. <i>New in version 0.19.2.</i>
<b>right_by</b> : <i>column name</i>	Field names to match on in the right DataFrame. <i>New in version 0.19.2.</i>
<b>suffixes</b> : <i>2-length sequence (tuple, list, ...)</i>	Suffix to apply to overlapping column names in the left and right side, respectively.
<b>tolerance</b> : <i>integer or Timedelta, optional, default None</i>	Select asof tolerance within this range; must be compatible with the merge index.

**allow\_exact\_matches** : boolean, default `True`

- If `True`, allow matching with the same ‘on’ value (i.e. less-than-or-equal-to / greater-than-or-equal-to)
- If `False`, don’t match the same ‘on’ value (i.e., strictly less-than / strictly greater-than)

**direction** : ‘backward’ (default), ‘forward’, or ‘nearest’

Whether to search for prior, subsequent, or closest matches.

*New in version 0.20.0.*

**Returns:**

**merged** : *DataFrame*

**See also:** `merge`, `merge_ordered`

```
1#!/usr/bin/env python3
2# -*- coding: utf-8 -*-
3"""
4Created on Wed Dec 20 19:00:09 2017
5
6@author: DemiWu
7"""
8
9import pandas as pd
10trades = pd.DataFrame({'time': pd.to_datetime(['20160525 13:30:00.023', '20160525 13:30:00.038', '20160525 13:30:00.048',
11                                              '20160525 13:30:00.048', '20160525 13:30:00.048']),
12                         'ticker': ['MSFT', 'MSFT', 'GOOG', 'GOOG', 'AAPL'],
13                         'price': [51.95, 51.95, 720.77, 720.92, 98.00],
14                         'quantity': [75, 155, 100, 100, 100]},
15                         columns=['time', 'ticker', 'price', 'quantity'])
16'''trades looks like
17          time ticker  price  quantity
18 0 2016-05-25 13:30:00.023    MSFT   51.95       75
19 1 2016-05-25 13:30:00.038    MSFT   51.95      155
20 2 2016-05-25 13:30:00.048    GOOG  720.77      100
21 3 2016-05-25 13:30:00.048    GOOG  720.92      100
22 4 2016-05-25 13:30:00.048    AAPL   98.00      100
23'''
24quotes = pd.DataFrame({'time': pd.to_datetime(['20160525 13:30:00.023', '20160525 13:30:00.023', '20160525 13:30:00.030',
25                                              '20160525 13:30:00.041', '20160525 13:30:00.048', '20160525 13:30:00.049',
26                                              '20160525 13:30:00.072', '20160525 13:30:00.075']),
27                         'ticker': ['GOOG', 'MSFT', 'MSFT', 'MSFT', 'GOOG', 'AAPL', 'GOOG', 'MSFT'],
28                         'bid': [720.50, 51.95, 51.97, 51.99, 720.50, 97.99, 720.50, 52.01],
29                         'ask': [720.93, 51.96, 51.98, 52.00, 720.93, 98.01, 720.88, 52.03]},
30                         columns=['time', 'ticker', 'bid', 'ask'])
31'''quotes looks like
32          time ticker    bid    ask
33 0 2016-05-25 13:30:00.023    GOOG  720.50  720.93
34 1 2016-05-25 13:30:00.023    MSFT   51.95   51.96
35 2 2016-05-25 13:30:00.030    MSFT   51.97   51.98
36 3 2016-05-25 13:30:00.041    MSFT   51.99   52.00
37 4 2016-05-25 13:30:00.048    GOOG  720.50  720.93
38 5 2016-05-25 13:30:00.049    AAPL   97.99  98.01
39 6 2016-05-25 13:30:00.072    GOOG  720.50  720.88
40 7 2016-05-25 13:30:00.075    MSFT   52.01  52.03
41'''
```

```
42
43 pd.merge_asof(trades, quotes, on='time', by='ticker')
44 '''the result is
45          time ticker  price  quantity      bid      ask
46 0 2016-05-25 13:30:00.023   MSFT  51.95       75  51.95  51.96
47 1 2016-05-25 13:30:00.038   MSFT  51.95      155  51.97  51.98
48 2 2016-05-25 13:30:00.048   GOOG  720.77      100 720.50 720.93
49 3 2016-05-25 13:30:00.048   GOOG  720.92      100 720.50 720.93
50 4 2016-05-25 13:30:00.048   AAPL  98.00      100      NaN      NaN
51 '''
52 #for the row with the time 13:30:00.038, the merged dataframe uses the last row in "quotes" which matches with
53 #the ticker "MSFT" and whose time is less than 13:30:00.038 (the one with the time 13:30:00.030)
54
55 pd.merge_asof(trades, quotes, on='time', by='ticker', tolerance=pd.Timedelta('2ms'))
56 '''the result is
57          time ticker  price  quantity      bid      ask
58 0 2016-05-25 13:30:00.023   MSFT  51.95       75  51.95  51.96
59 1 2016-05-25 13:30:00.038   MSFT  51.95      155      NaN      NaN
60 2 2016-05-25 13:30:00.048   GOOG  720.77      100 720.50 720.93
61 3 2016-05-25 13:30:00.048   GOOG  720.92      100 720.50 720.93
62 4 2016-05-25 13:30:00.048   AAPL  98.00      100      NaN      NaN
63 '''
64 #because 13:30:00.030 is away from 13:30:00.038 by more than 2ms, the merged does not uses the corresponding data
65 #from "quotes"
66
```

```
00
67 pd.merge_asof(trades, quotes, on='time', by='ticker', tolerance=pd.Timedelta('10ms'), allow_exact_matches=False)
68 '''the result is
69
70      time ticker  price  quantity    bid    ask
71 0 2016-05-25 13:30:00.023   MSFT  51.95       75    NaN    NaN
72 1 2016-05-25 13:30:00.038   MSFT  51.95      155  51.97  51.98
73 2 2016-05-25 13:30:00.048   GOOG  720.77      100    NaN    NaN
74 3 2016-05-25 13:30:00.048   GOOG  720.92      100    NaN    NaN
75 4 2016-05-25 13:30:00.048   AAPL   98.00      100    NaN    NaN
76 '''
76 #only the row with the time 13:30:00.038 has bid-ask data from "quotes" because it uses data from near rows rather than
77 #exact match
78
79 pd.merge_asof(trades, quotes, on='time', by='ticker', direction='nearest')
80 '''the result is
81
82      time ticker  price  quantity    bid    ask
83 0 2016-05-25 13:30:00.023   MSFT  51.95       75  51.95  51.96
84 1 2016-05-25 13:30:00.038   MSFT  51.95      155  51.99  52.00
85 2 2016-05-25 13:30:00.048   GOOG  720.77      100  720.50  720.93
86 3 2016-05-25 13:30:00.048   GOOG  720.92      100  720.50  720.93
87 4 2016-05-25 13:30:00.048   AAPL   98.00      100  97.99  98.01
88 '''
88 #the row with the time 13:30:00.038 now uses the bid-ask data from "quotes" in the row with time 13:00:00.041, \
89 #which is the closest to 13:30:00.038
90
91 #the other arguments are the same as in pandas.merge function
```

<b>Function Type</b>	<b># Functions</b>
General functions	26
Series	283
DataFrame	197
Panel	107
Index	69
Other indices	87
Window	35
Groupby	63
Resampling	29
Style	16
General utility functions	58
Total	970

# Course Outline

- Data management with Python (3 weeks)
  - Packages
    - Numpy/Scipy
    - Dataframe
    - Pandas
  - OOP
    - Class
  - Geolocation
  - Multiprocessing
  - Pickle, Json, HDF5 (.h5)
  - Web scraping (regular expression, selenium)

# Object

- An Object is a bit of self-contained Code and Data
- A key aspect of the Object approach is to break the problem into smaller understandable parts (divide and conquer)
- Objects have boundaries that allow us to ignore un-needed detail
- We have been using objects all along: String Objects, Integer Objects, Dictionary Objects, List Objects...

# Definitions

- **Class** - a template - Dog
- **Method or Message** - A defined capability of a class - bark()
- **Field or attribute**- A bit of data in a class - length
- **Object or Instance** - A particular instance of a class - Snoopy

# Terminology: Class

Defines the abstract characteristics of a thing (object), including:

1. the thing's characteristics (its attributes, **fields** or **properties**)
2. the thing's behaviors (the things it can do, or **methods**, operations or features).

One might say that a **class** is a **blueprint** or factory that describes the nature of something.

For example, the **class** `Dog` would consist of traits shared by all dogs, such as breed and fur color (characteristics), and the ability to bark and sit (behaviors).

# Terminology: Class



A **pattern** (exemplar) of a class. The class of Dog defines all possible dogs by listing the characteristics and behaviors they can have; the object Snoopy is one particular dog, with particular versions of the characteristics. A Dog has fur; Snoopy has white fur.

# Terminology: Instance



One can have an **instance** of a class or a particular object. The **instance** is the actual object created at runtime. The Snoopy object is an **instance** of the Dog class. The set of values of the attributes of a particular object is called its **state**. The object consists of state and the behavior that's defined in the object's class.

Object and Instance are often used interchangeably.

# Terminology: Method



An object's abilities. In language, **methods** are verbs. Snoopy, being a Dog, has the ability to bark. So bark() is one of Snoopy's methods. He may have other **methods** as well, for example sit() or eat() or walk(). Within the program, using a **method** usually affects only one particular object; all Dogs can bark, but you need only one particular dog to do the barking

Method and Message are often used interchangeably.

# A Sample Class

class is a reserved word.

Each PartyAnimal object has a bit of code.

Tell the object to run the party() code.

```
class PartyAnimal:  
    x = 0  
  
    def party(self) :  
        self.x = self.x + 1  
        print("So far",self.x)  
  
an = PartyAnimal()  
  
an.party()  
an.party()  
an.party()
```

This is the template for making PartyAnimal objects.

Each PartyAnimal object has a bit of data.

Create a PartyAnimal object.

PartyAnimal party an

run party() \*within\* the object an

# Definitions



- **Class** - a template - Dog
- **Method or Message** - A defined capability of a class - bark()
- **Object or Instance** - A particular instance of a class - Snoopy

Playing with dir() and type()

# A Nerdy Way to Find Capabilities

- The `dir()` command lists capabilities
- Ignore the ones with underscores - these are used by Python itself
- The rest are real operations that the object can perform
- It is like `type()` - it tells us something \*about\* a variable

```
>>> x = list()
>>> type(x)
<type 'list'>
>>> dir(x)
['__add__', '__class__',
 '__contains__', '__delattr__',
 '__delitem__', '__delslice__',
 '__doc__', '__eq__',
 '__getitem__', '__setslice__',
 '__str__', 'append', 'count',
 'extend', 'index', 'insert', 'pop',
 'remove', 'reverse', 'sort'
```

# Try dir() with a String

```
>>> y = "Hello there"  
>>> dir(y)  
['__add__', '__class__', '__contains__', '__delattr__', '__doc__',  
'__eq__', '__ge__', '__getattribute__', '__getitem__',  
'__getnewargs__', '__getslice__', '__gt__', '__hash__',  
'__init__', '__le__', '__len__', '__lt__', '__repr__', '__rmod__',  
'__rmul__', '__setattr__', '__str__', 'capitalize', 'center', 'count',  
'decode', 'encode', 'endswith', 'expandtabs', 'find', 'index',  
'isalnum', 'isalpha', 'isdigit', 'islower', 'isspace', 'istitle', 'isupper',  
'join', 'ljust', 'lower', 'lstrip', 'partition', 'replace', 'rfind', 'rindex',  
'rjust', 'rpartition', 'rsplit', 'rstrip', 'split', 'splitlines', 'startswith',  
'strip', 'swapcase', 'title', 'translate', 'upper', 'zfill']
```

```
class PartyAnimal:  
    x = 0  
  
    def party(self):  
        self.x = self.x + 1  
        print ("So far",self.x)  
  
an = PartyAnimal()  
  
print ("Type", type(an))  
print ("Dir ", dir(an))
```

We can use `dir()` to find the “capabilities” of *our* newly created class.

```
$ python party2.py  
Type <type 'instance'>  
Dir ['__doc__', '__module__',  
'party', 'x']
```

# Constructor

- The primary purpose of the constructor is to set up some instance variables to have the proper initial values when the object is created

```
class PartyAnimal:  
    x = 0  
  
    def __init__(self):  
        print ("I am constructed")  
  
    def party(self) :  
        self.x = self.x + 1  
        print ("So far",self.x)  
  
    def __del__(self):  
        print ("I am destructured", self.x)  
  
an = PartyAnimal()  
an.party()  
an.party()  
an.party()
```

```
$ python party2.py  
I am constructed  
So far 1  
So far 2  
So far 3  
I am destructured 3
```

The constructor and destructor are optional. The constructor is typically used to set up variables. The destructor is seldom used.

# Constructor



- In **object-oriented programming**, a **constructor** in a class is a special block of statements called when an **object is created**

# Many Instances

- We can create **lots of objects** - the class is the template for the object
- We can store each **distinct object** in its own variable
- We call this having multiple **instances** of the same class
- Each **instance** has its own copy of the **instance variables**

```
class PartyAnimal:  
    x = 0  
    name = ""  
    def __init__(self, nam):  
        self.name = nam  
        print (self.name,"constructed")
```

```
def party(self) :  
    self.x = self.x + 1  
    print (self.name,"party count",self.x)
```

```
s = PartyAnimal("Sally")  
s.party()
```

```
j = PartyAnimal("Jim")  
j.party()  
s.party()
```

Constructors can have additional parameters. These can be used to setup instance variables for the particular instance of the class (i.e. for the particular object).

# Definitions



- **Class** - a template - Dog
- **Method or Message** - A defined capability of a class - bark()
- **Object or Instance** - A particular instance of a class - Lassie
- **Constructor** - A method which is called when the instance / object is created

# Inheritance

<https://docs.python.org/3/tutorial/classes.html#inheritance>

<http://www.ibiblio.org/g2swap/bytewofpython/read/inheritance.html>

# Inheritance

- When we make a new class - we can reuse an existing class and **inherit** all the capabilities of an existing class and then add our own little bit to make our new class
- Another form of store and reuse
- Write once - reuse many times
- The new class (child) has all the capabilities of the old class (parent) - and then some more

# Terminology: Inheritance

‘Subclasses’ are more specialized versions of a class, which **inherit** attributes and behaviors from their parent classes, and can introduce their own.

[http://en.wikipedia.org/wiki/Object-oriented\\_programming](http://en.wikipedia.org/wiki/Object-oriented_programming)

```
class PartyAnimal:  
    x = 0  
    name = ""  
  
    def __init__(self, nam):  
        self.name = nam  
        print (self.name,"constructed")
```

```
def party(self) :  
    self.x = self.x + 1  
    print (self.name,"party count",self.x)
```

```
class FootballFan(PartyAnimal):  
    points = 0  
  
    def touchdown(self):  
        self.points = self.points + 7  
        self.party()  
        print (self.name,"points",self.points)
```

```
s = PartyAnimal("Sally")  
s.party()  
  
j = FootballFan("Jim")  
j.party()  
j.touchdown()
```

**FootballFan** extends **PartyAnimal**. It has all the capabilities of **PartyAnimal** and more.

# SUMMARY



# Classes

```
class MyVector: """A simple vector class."""

    num_created = 0

    def __init__(self,x=0,y=0):
        self.__x = x
        self.__y = y
        MyVector.num_created += 1

    def get_size(self):
        return self.__x+self.__y

    @staticmethod
    def get_num_created
        return MyVector.num_created
```

## #USAGE OF CLASS MyVector

```
print (MyVector.num_created)
v = MyVector()
w = MyVector(0.23,0.98)
print (w.get_size())
bool = isinstance(v, MyVector)
```

Output:

0

1.21



# I/O

```
import os
print (os.getcwd()) #get "."
os.chdir('..')
import glob # file globbing
lst = glob.glob('*.*txt') # get list of files
import shutil # mngmt tasks
shutil.copyfile('a.py','a.bak')

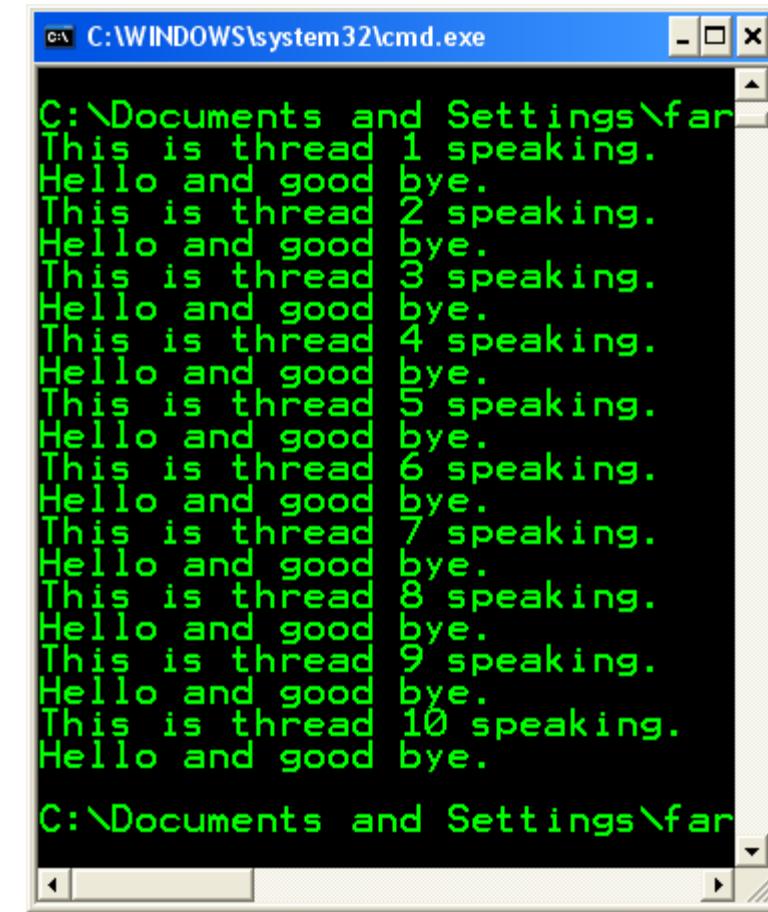
import pickle # serialization logic
ages = {"ron":18,"ted":21}
pickle.dump(ages,fout)
# serialize the map into a writable file
ages = pickle.load(fin)
# deserialize map from areadable file

# read binary records from a file
from struct import *
fin = None
try:
    fin = open("input.bin","rb")
    s = f.read(8)#easy to read in
    while (len(s) == 8):
        x,y,z = unpack(">HH<L", s)
        print ("Read record: " \
               "%04x %04x %08x"%(x,y,z))
        s = f.read(8)
except IOError:
    pass
if fin: fin.close()
```



# Threading in Python

```
import threading  
theVar = 1  
  
class MyThread ( threading.Thread ):  
    def run ( self ):  
        global theVar  
        print('This is thread ' + \  
              str ( theVar ) + ' speaking.')  
        print('Hello and good bye.' )  
        theVar = theVar + 1  
  
for x in range ( 10 ):  
    MyThread().start()
```



```
C:\Documents and Settings\far  
This is thread 1 speaking.  
Hello and good bye.  
This is thread 2 speaking.  
Hello and good bye.  
This is thread 3 speaking.  
Hello and good bye.  
This is thread 4 speaking.  
Hello and good bye.  
This is thread 5 speaking.  
Hello and good bye.  
This is thread 6 speaking.  
Hello and good bye.  
This is thread 7 speaking.  
Hello and good bye.  
This is thread 8 speaking.  
Hello and good bye.  
This is thread 9 speaking.  
Hello and good bye.  
This is thread 10 speaking.  
Hello and good bye.  
C:\Documents and Settings\far
```

# Graphics to explore the data

Seaborn package is built on matplotlib but provides high level interface for drawing attractive statistical graphics, similar to ggplot2 library in R. It specifically targets statistical data visualization

To show graphs within Python notebook include inline directive:

```
In [ ]: %matplotlib inline
```

# Graphics

description	
distplot	histogram
barplot	estimate of central tendency for a numeric variable
violinplot	similar to boxplot, also shows the probability density of the data
jointplot	Scatterplot
regplot	Regression plot
pairplot	Pairplot
boxplot	boxplot
swarmplot	categorical scatterplot
factorplot	General categorical plot

# Basic statistical Analysis

statsmodel and scikit-learn - both have a number of function for statistical analysis

The first one is mostly used for regular analysis using R style formulas, while scikit-learn is more tailored for Machine Learning.

statsmodels:

- linear regressions
- ANOVA tests
- hypothesis testings
- many more ...

scikit-learn:

- kmeans
- support vector machines
- random forests
- many more ...

See examples in the Tutorial Notebook