# Financial Bigdata and Python

## 3. NumPy, Matplotlib

# Jump to Python
https://wikidocs.net/book/1



Pahk Eungyong (pahkey@gmail.com)
An introduction to Python published on Wikidocs (https://wikidocs.net)

A handy tool for running simple scripts.
Just type 'ipython3' to begin the interactive shell.
To exit, press the keys Ctrl and D keys together.

Anaconda Powershell

Run Jupyter Notebook

If the script you want to run isn't in one sentence,
you're better off using a Jupyter Notebook rather than Anaconda PowerShell.

| Anaconda PowerShell | Description |
| --- | --- |
| abs(3)<br>abs(-3)<br>abs(-1.2) | returns the absolute value of an input |
| pow(2, 4)<br>pow(3, 3) | returns the power of an input<br>pow(a,b) → a^b |
| round(4.6)<br>round(4.2)<br>round(5.678, 2) | rounding function<br><br>specifies the number of decimal places to be rounded off. |
| max([1, 2, 3])<br>min([1, 2, 3])<br>sum([1, 2, 3]) | returns the maximum of the list<br>returns the minimum of the list<br>returns the sum of the list |
| sorted([3, 1, 2])<br>sorted(['a', 'c', 'b']) | sorts the input list |

Exercise 1)

Let's find the sum, the maximum and minimum values of the following list.
[-8, 2, 7, 5, -3, 5, 0, 1]

Exercise 2)

Here is the result when 17/3 is output in Python:

>>> 17 / 3
5.666666666666667

Let's display the result 5.666666666666667 rounded up to 4 decimal places.

| Anaconda PowerShell | Description |
|---|---|
| import numpy as np | imports the package NumPy<br>'np' is usually used as the abbreviation of NumPy. |
| a = np.array([1,2,3,4])<br>b = np.array([5,6,7,8]) | create arrays a and b of size 4 |
| c = a+b<br>d = a-b<br>e = a*b<br>f = a/b | Arithmetic operations between two matrices a and b of the same size<br>(calculated elementwise between elements in the same position) |
| g = a+2<br>h = a-2<br>i= 2*a<br>j =a/2 | Arithmetic operations between an array and a scalar<br>(Such an operation applies to all elements of the matrix) |
| How to implement matrix addition in python without NumPy?<br><br>c = []<br>for e_a, e_b in zip(a,b) :<br>    e_c = e_a + e_b<br>    c.append(e_c) | Very inefficient way.<br>In other words, it's really slow.<br><br>So if it is not for study or research purposes,<br>do not reinvent the wheel! |

| Anaconda PowerShell | Description |
|---|---|
| A = np.array([[1,1,1],[1,1,1],[1,1,1],[1,1,1]])<br>A = np.ones([4,3]) | creates a matrix of size 4x3<br>creates a matrix of size 4x3 where all elements are 1. |
| z = np.array([5,6])<br>A+z | creates a matrix of size 2<br>An error occurs because the matrices have different sizes. |
| a = np.array([1,2,3])<br>A+a | It is expected that an error occurs because the size of the matrix is different, but this code runs well in the way of sequentially adding the elements of 'a' to each row of 'A'.<br><br>Even if it's annoying that a code doesn't run because of an error, it's not severe. However, it is rather dangerous if it is not performing as expected, but there is no error. |

- Useful array creation functions to know (array = matrix + vector + tensor)

zeros([n,m]): creates a matrix of size n×m where all elements are 0
ones([n,m]): creates a matrix of size n×m where all elements are 1.
arange(a,b): similar to the range function.
        ex) arange(0,5) → array([0,1,2,3,4])
linspace(a,b,n): creates a matrix where n evenly spaced numbers from 'a' to 'b'.
        ex) linspace(0,2,5) → array([0,0.5,1,1.5,2])

| Anaconda PowerShell | Description |
|---|---|
| a = np.array([1,2,3,4])<br>a.shape<br>a.ndim | the shape of a : (4,)<br>the dimension of a : 1 |
| b = np.array([[1,2,3,4],[5,6,7,8]])<br>b.shape<br>b.ndim | the shape of b : (2,4)<br>the dimension of b : 2 |
| c = np.array([[1,2],[3,4]])<br>c.shape<br>c.ndim | the shape of c : (2,2)<br>the dimension of c : 2 |
| d = a.reshape([2,2])<br>e = b.reshape(8) | changes the shape of a from (4,) to (2,2)<br>changes the shape of b from (2,4) to (8,) |
| f = b.reshape(-1)<br>g = a.reshape([2,-1])<br>h = a.reshape([-1,4]) | changes the shape of b from (2,4) to (8,)<br>changes the shape of a from (4,) to (2,2)<br>changes the shape of a from (4,) to (1,4) |
| np.vstack([a,b])<br>np.hstack([b,c]) | stacks matrices with the same number of columns top and bottom<br>stacks a matrix with the same number of rows on both sides |

| Anaconda PowerShell | Description |
|---|---|
| a>2 | the logical matrix generated by whether the element of 'a' is greater than 2 |
| a[a>2] | the elements of 'a' greater than 2 |
| | |
| b>3 | |
| b[b>3] | |
| | |
| a[0] | accesses the 0th element of a |
| a[:3] | accesses 0-2 elements of a |
| a[-1] | accesses the last element of a |
| | |
| b[0,0] | accesses the 0th row, the 0th column element of b |
| b[0,:2] | accesses the 0th row, the 0-1 column of b |
| b[:,1:3] | accesses all the rows, the 1-2 columns of b |
| | |
| x = np.linspace(0,1,5) | creates a vector by dividing [0,1] into 4 equal parts |
| y = np.linspace(1,3,6) | creates a vector by dividing [1,3] into 5 equal parts |
| | |
| X, Y = np.meshgrid(x,y) | creates a 2D mesh based on x,y vectors |
| | |
| X | X is the x-coordinate of the mesh |
| Y | Y is the y-coordinate of the mesh |

| Anaconda PowerShell | Description |
|---|---|
| a = np.array([1,2,3,4])<br><br>len(a)<br>np.sum(a)<br>np.mean(a)<br>np.var(a)<br>np.std(a)<br>np.max(a)<br>np.min(a)<br>np.argmax(a)<br>np.argmin(a)<br>np.median(a)<br>np.percentile(a,percentile) | <br><br>the number of elements in a<br>the sum of all elements of a<br>the average value of a<br>the variance of a<br>the deviation value of a<br>the maximum value of a<br>the minimum value of a<br>the maximum position of a<br>the position of the minimum value of a<br>the median of a<br>the percentile of a. (ex) percentil(a,25) : 1st quartile |

There are so many NumPy functions. It is impossible and undesirable to cover them all in this class.
(Like not memorizing an English dictionary during an English class..).

NumPy official user guide:
https://numpy.org/devdocs/user/index.html

Data Science School (by Kim Do-hyung) :
https://datascienceschool.net

Exercise 3)

Solve the problems for the following matrix.

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

a) Index the value 7.
b) Index the value 14.
c) Slice the array [6, 7].
d) Slice the array [7, 12].
e) Slice the array [[3, 4], [8, 9]].

Exercise 4)

Create the following array using the functions introduced so far.

array([[    0.,    0.,    0.,   1.,     1.],
       [    0.,    0.,    0.,   1.,     1.],
       [    0.,    0.,    0.,   1.,     1.],
       [  10.,   20.,   30.,   40.,    50.],
       [  60.,   70.,   80.,   90.,   100.],
       [ 110.,  120.,  130.,  140.,   150.],
       [    0.,    0.,    0.,   1.,     1.],
       [    0.,    0.,    0.,   1.,     1.],
       [    0.,    0.,    0.,   1.,     1.],
       [  10.,   20.,   30.,   40.,    50.],
       [  60.,   70.,   80.,   90.,   100.],
       [ 110.,  120.,  130.,  140.,   150.]]).

Exercise 5)

Solve the problems for the following matrix.

x = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10,
           11, 12, 13, 14, 15, 16, 17, 18, 19, 20])

a) Find multiples of 3.
b) Find the numbers that when divided by 4 leaves 1 left.
c) Find the numbers among multiples of 3 that when divided by 4 leaves 1 left.

| Anaconda PowerShell | Description |
|---|---|
| import numpy as np | |
| a = np.array( [ 1,2,3,4 ] )<br>a | creates a vector |
| b = np.array( [[ 1,2,3,4 ]] )<br>b | creates a row vector (= a matrix with row length 1) |
| c = np.array( [[ 1 ],[ 2 ],[ 3 ],[ 4 ]] )<br>c | creates a column vector (= a matrix of column length 1) |
| a.ndim<br>a.shape | ndim of scalars: 0<br>ndim of vectors: 1<br>ndim of matrices (column vector, row vector): 2 |
| b.ndim<br>b.shape | |
| c.ndim<br>c.shape | shape of scalars : ()<br>shape of vectors: (n,)<br>shape of matrices (column vector, row vector): (n,m) |

| Anaconda PowerShell | Description |
|---|---|
| | **\* Method 1** |
| a.reshape([4,1]) | converts 'vector' to 'column vector' using the reshape function |
| a.reshape([1,4]) | converts 'vector' to 'row vector' |
| | |
| b.reshape(4) | converts 'row vector' to 'vector' |
| b.reshape([4,1]) | converts 'row vector' to 'column vector' |
| | |
| c.reshape(4) | converts 'column vector' to 'vector' |
| c.reshape([1,4]) | converts 'column vector' to 'row vector' |
| | |
| | **\* Method 2** |
| a[:,np.newaxis] | converts 'column vector' to 'vector' by adding axis 1 |
| a[np.newaxis,:] | converts 'row vector' to 'vector' by adding axis 0 |
| | |
| b.flatten() | converts 'row vector' to 'vector' using flatten function |
| b.transpose() | converts 'row vector' to 'column vecotor' by transposing it |
| | |
| c.flatten() | converts 'column vector' to 'vector' |
| c.transpose() | converts 'column vector' to 'row vector' |

| Anaconda PowerShell | Description |
|---|---|
| sa = a.sum()<br>sa<br>type(sa)<br>sa.ndim, sa.shape | returns the sum of all elements of 'vector'<br><br>Note that the return value is of type 'np.int', not of type 'int'.<br>So you can still get the ndim and shape for the variable.<br>Also notice that ndim is reduced from 1 to 0. |
| sa = a.sum(keepdims=True)<br>sa<br>type(sa)<br>sa.ndim, sa.shape | sums all elements while preserving dimension.<br>'a' is a vector of the shape (4,), so 'sa' has the shape (1,)<br>(Note that () ≠ (1,)) |
| sb = b.sum()<br>sb<br>sb.ndim, sb.shape | |
| sb = b.sum(axis=1)<br>sb<br>sb.ndim, sb.shape | sums all elements along the 1st (specific) axis<br>'b' has the shape (1,4), so 'sb' has the shape (1,)<br>(the 1st axis has disappeared and the 0th axis has been left.)<br>Meanwhile, ndim is reduced from 2 to 1. |
| sb = b.sum(axis=1,keepdims=True)<br>sb<br>sb.ndim, sb.shape | sums all elements along the 1st axis while preserving dimension In that case, the shape of 'sb' will be (1,1), not rather (1,).<br>So ndim has been also kept at 2. |
| sb = b.sum(axis=0)<br>sb<br>sb.ndim, sb.shape | For a 'column vector' like 'b', summing the elements along the 0th axis changes nothing other than changing 'ndim'. |

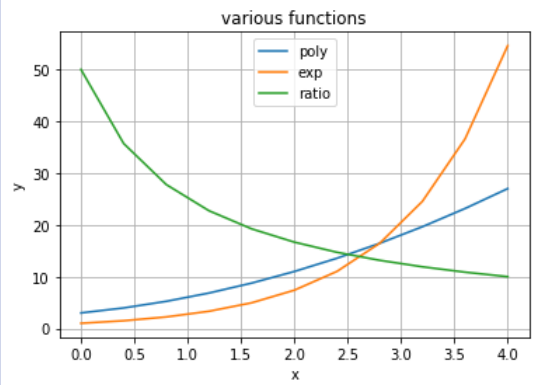| Anaconda PowerShell | Description |
|---|---|
| A = np.array([[1,2,3],[4,5,6]])<br>A<br>A.ndim, A.shape<br>A.flatten()<br>A.transpose() | matrix of the shape (2,3)<br><br><br>flattens the matrix by row-major order<br>swaps the 0th axis and 1st axis |
| A[np.newaxis,:,:].shape<br>A[:,np.newaxis,:].shape<br>A[:,:,np.newaxis].shape | A new axis is added as the 0th axis.<br>A new axis is added as the 1st axis.<br>A new axis is added as the 2nd axis. |
| sb = A.sum(axis=0)<br>sb<br>sb.ndim, sb.shape | sums all elements along the 0th axis |
| sb = A.sum(axis=1)<br>sb<br>sb.ndim, sb.shape | sums all elements along the 1st axis |
| sb = A.sum(axis=1,keepdims=True)<br>sb<br>sb.ndim, sb.shape | sums all elements along the 1st axis while preserving dimension. |

| Anaconda PowerShell | Description |
|---|---|
| a = np.array([1,2,3,4])<br>a.shape<br>a.ndim | shape: (4,)<br>ndim: 1 |
| b = np.array([[1,2,3,4],[5,6,7,8]])<br>b.shape<br>b.ndim | shape: (2,4)<br>ndim: 2 |
| c = np.array([[1,2],[3,4]])<br>c.shape<br>c.ndim | shape: (2,2)<br>ndim: 2 |
| np.vstack([a,b])<br>np.hstack([b,c]) | stacks matrices with the same number of columns top and bottom<br>stacks matrices with the same number of rows on both sides |
| np.concatenate([a,b],axis=0) | a: (4,) / b: (2,4)<br>    : stacking the matrices on axis 0<br>    (An error occur because the dimensions of a and b are different) |
| a2 = a[np.newaxis,:]<br>np.concatenate([a2,b],axis=0) | a: (1,4) / b: (2,4) ( = np.hstack([a,b]) )<br>    : stacking the matrices on axis 0 |
| np.concatenate([b,c],axis=1) | b: (2,4) / c: (2,2) ( = np.vstack([b,c]) )<br>    : stacking the matrices on axis 1 |

| Anaconda PowerShell | Description |
|---|---|
| x = np.arange(4)<br>print('x',x.shape); print(x) | [0,1,2,3]    (4,) |
| xx = x.reshape(4,1)<br>print('xx',xx.shape); print(xx) | [[0],          (4,1)<br> [1],<br> [2],<br> [3]] |
| y = np.ones(5)<br>print('y',y.shape); print(y) | [1,1,1,1,1]   (5,) |
| z = np.ones((3,4))<br>print('z',z.shape); print(z) | [[1,1,1,1],    (3,4)<br> [1,1,1,1],<br> [1,1,1,1]] |
| x+y | (4,) + (5,)  Error! |
| | Broadcasting is applied to below examples |
| xx+y<br>(xx+y).shape | (4,1) + (5,)        = (4,1) + (1,5)          = (4,5) + (4,5)<br>[[0], + [1,1,1,1,1] = [[0], + [[1,1,1,1,1]] = [[0,0,0,0,0], + [[1,1,1,1,1],<br> [1],                    [1],                          [1,1,1,1,1],      [1,1,1,1,1],<br> [2],                    [2],                          [2,2,2,2,2],      [1,1,1,1,1],<br> [3]]                    [3]]                          [3,3,3,3,3]]     [1,1,1,1,1]] |
| x+z<br>(x+z).shape | (4,) + (3,4)            = (1,4) + (3,4)             = (3,4) + (3,4)<br>[0,1,2,3] + [[1,1,1,1],  = [[0,1,2,3]] + [[1,1,1,1],   = [[0,1,2,3], + [[1,1,1,1],<br>              [1,1,1,1],                      [1,1,1,1],      [0,1,2,3],   [1,1,1,1],<br>              [1,1,1,1]]                      [1,1,1,1]]     [0,1,2,3]]   [1,1,1,1]] |

Exercise 6) Dimension reduction operations (sum, max, etc.)

Create a  5 x 6 matrix composed of real numbers and obtain the following values for the matrix.

1) The maximum value of the whole
2) Sum of each row
3) Maximum value in each row
4) Average of each column
5) Minimum in each column

| Anaconda PowerShell | Description |
|---|---|
| import numpy as np<br>import matplotlib.pylab as plt<br><br>xarr = np.linspace(0,4,11)<br>yarr1 = xarr**2 + 2*xarr + 3<br>yarr2 = np.exp(xarr)<br>yarr3 = 50/(1+xarr) | MatplotLib : popular package for drawing using Python<br>pylab : introduces MATLAB interface to MatplotLib |
| plt.figure(figsize=(6,4))<br>plt.plot(xarr,yarr1)<br>plt.plot(xarr,yarr2)<br>plt.plot(xarr,yarr3)<br>plt.xlabel('x')<br>plt.ylabel('y')<br>plt.title('various functions')<br>plt.legend(['poly','exp','ratio'])<br>plt.grid()<br>plt.show() | creates a figure window with size (6,4)<br>plots the graph of (xarr,yarr1)<br>plots the graph of (xarr,yarr2)<br>plots the graph of (xarr,yarr3)<br>x-axis labeling<br>y-axis labeling<br>titles the graphs<br>makes legends<br>makes the grid appear<br>shows the graphs<br> |

| Anaconda PowerShell | Description |
|---|---|

```
plt.figure(figsize=(6,4))

plt.plot(xarr,yarr1,'-o')
plt.plot(xarr,yarr2,'--X',ms=10)
plt.plot(xarr,yarr3,':D',lw=2)

plt.xlabel('x')
plt.ylabel('y')
plt.title('various functions')
plt.legend(['poly','exp','ratio'])
plt.grid()
plt.show()
```
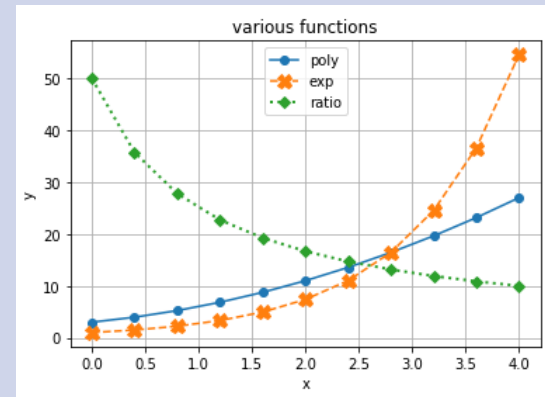
※ line style

- solid line
-- dashed line
-. dash-dot line
: dotted lined

※ plot function option

| Option | Abbreviation | Description |
|---|---|---|
| color | c | line color |
| linewidth | lw | line width |
| linestyle | ls | line style |
| marker | | marker type |
| markersize | ms | marker size |
| markeredgecolor | mec | marker line color |
| markeredgewidth | mew | marker line width |
| markerfacecolor | mfc | marker inner color |

**※ color**  http://Matplotlib.org/examples/color/named_colors.html

| color | abbreviation |
|-------|--------------|
| blue | b |
| green | g |
| red | r |
| cyan | c |
| magenta | m |
| yellow | y |
| black | k |
| white | w |

**※ marker**  https://matplotlib.org/3.2.1/api/markers_api.html

| marker type | |
|-------------|---|
| o | circle |
| ^ | triangle |
| s | square |
| P | cross |
| * | star |
| X | X |
| D | diamond |
| $...$ (ex. $f$) | desired character |

| Anaconda PowerShell | Description |
|---|---|
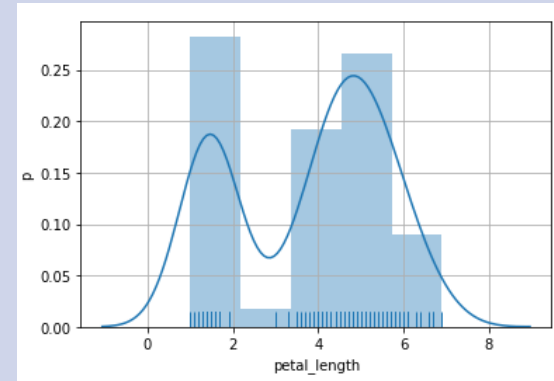| ```
import matplotlib.pylab as plt
import seaborn as sns

iris = sns.load_dataset('iris')
iris
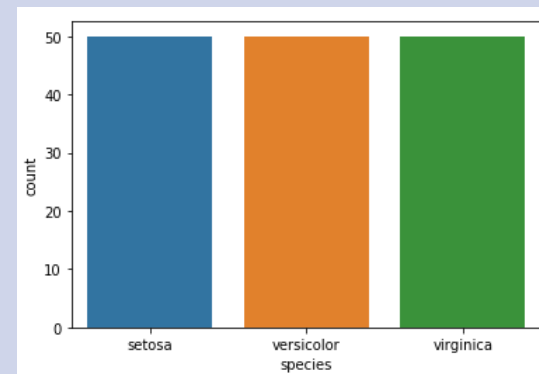

fig=plt.figure(figsize=(6,4))
sns.displot(data=iris,
        x='petal_length',
        kde=True, rug=True)
plt.ylabel('p')
plt.grid()
plt.show()

fig=plt.figure(figsize=(6,4))
sns.countplot(data=iris,
        x='species')
plt.show()
``` | seaborn: an extension of MatplotLib for statistical visualization<br><br><br><br>displot() function is useful for drawing histogram, density, and rugs.<br><br><br><br>countplot() function visualizes how many data each class has |