

# Python For Data Science Cheat Sheet

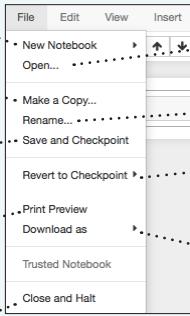
## Jupyter Notebook

Learn More Python for Data Science Interactively at [www.DataCamp.com](http://www.DataCamp.com)



### Saving/Loading Notebooks

Create new notebook



Make a copy of the current notebook

Save current notebook and record checkpoint

Preview of the printed notebook

Close notebook & stop running any scripts

Open an existing notebook

Rename notebook

Revert notebook to a previous checkpoint

Download notebook as

- IPython notebook
- Python
- HTML
- Markdown
- reST
- LaTeX
- PDF

### Writing Code And Text

Code and text are encapsulated by 3 basic cell types: markdown cells, code cells, and raw NBConvert cells.

#### Edit Cells

Cut currently selected cells to clipboard

Paste cells from clipboard above current cell

Paste cells from clipboard on top of current cell

Revert "Delete Cells" invocation

Merge current cell with the one above

Move current cell up

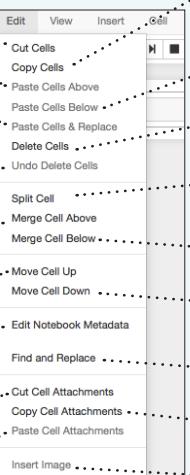
Adjust metadata underlying the current notebook

Remove cell attachments

Paste attachments of current cell

#### Insert Cells

Add new cell above the current one



Copy cells from clipboard to current cursor position

Paste cells from clipboard below current cell

Delete current cells

Split up a cell from current cursor position

Merge current cell with the one below

Move current cell down

Find and replace in selected cells

Copy attachments of current cell

Insert image in selected cells

### Working with Different Programming Languages

Kernels provide computation and communication with front-end interfaces like the notebooks. There are three main kernels:



IPython



IRkernel



Julia

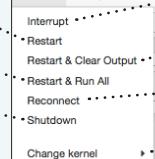
Installing Jupyter Notebook will automatically install the IPython kernel.

Restart kernel

Restart kernel & run all cells

Restart kernel & run all cells

Kernel Widgets Help



Interrupt kernel

Interrupt kernel & clear all output

Connect back to a remote notebook

Run other installed kernels

### Command Mode:



### Edit Mode:



### Executing Cells

Run selected cell(s)

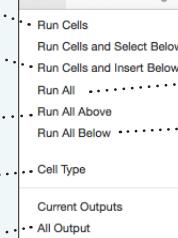
Run current cells down and create a new one above

Run all cells above the current cell

Change the cell type of current cell

toggle, toggle scrolling and clear all output

Cell Kernel Widgets



Run current cells down and create a new one below

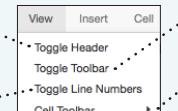
Run all cells

Run all cells below the current cell  
toggle, toggle scrolling and clear current outputs

### View Cells

Toggle display of Jupyter logo and filename

Toggle line numbers in cells



Toggle display of toolbar

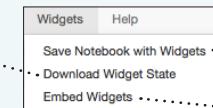
Toggle display of cell action icons:  
- None  
- Edit metadata  
- Raw cell format  
- Slideshow  
- Attachments  
- Tags

### Widgets

Notebook widgets provide the ability to visualize and control changes in your data, often as a control like a slider, textbox, etc.

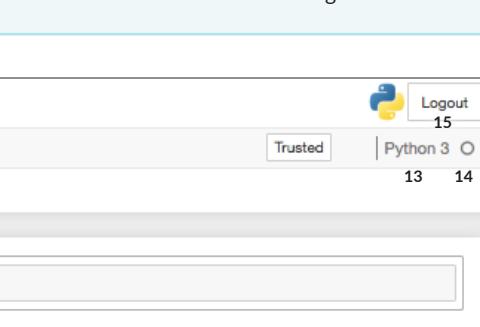
You can use them to build interactive GUIs for your notebooks or to synchronize stateful and stateless information between Python and JavaScript.

Download serialized state of all widget models in use



Save notebook with interactive widgets

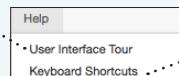
Embed current widgets



1. Save and checkpoint
2. Insert cell below
3. Cut cell
4. Copy cell(s)
5. Paste cell(s) below
6. Move cell up
7. Move cell down
8. Run current cell
9. Interrupt kernel
10. Restart kernel
11. Display characteristics
12. Open command palette
13. Current kernel
14. Kernel status
15. Log out from notebook server

### Asking For Help

Walk through a UI tour



List of built-in keyboard shortcuts

Edit the built-in keyboard shortcuts



Notebook help topics

Description of markdown available in notebook



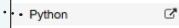
Information on unofficial Jupyter Notebook extensions

Python help topics



IPython help topics

NumPy help topics



SciPy help topics

Matplotlib help topics



SymPy help topics

Pandas help topics



About Jupyter Notebook



### Insert Cells



Add new cell below the current one



Learn Python for Data Science Interactively





# Python For Data Science Cheat Sheet

## Importing Data

Learn Python for data science interactively at [www.DataCamp.com](http://www.DataCamp.com)



### Importing Data in Python

Most of the time, you'll use either NumPy or pandas to import your data:

```
>>> import numpy as np  
>>> import pandas as pd
```

### Help

```
>>> np.info(np.ndarray.dtype)  
>>> help(pd.read_csv)
```

### Text Files

#### Plain Text Files

```
>>> filename = 'huck_finn.txt'  
>>> file = open(filename, mode='r')  
>>> text = file.read()  
>>> print(file.closed)  
>>> file.close()  
>>> print(text)
```

Open the file for reading  
Read a file's contents  
Check whether file is closed  
Close file

#### Using the context manager with

```
>>> with open('huck_finn.txt', 'r') as file:  
    print(file.readline())  
    print(file.readline())  
    print(file.readline())
```

Read a single line

### Table Data: Flat Files

#### Importing Flat Files with numpy

##### Files with one data type

```
>>> filename = 'mnist.txt'  
>>> data = np.loadtxt(filename,  
                    delimiter=',',  
                    skiprows=2,  
                    usecols=[0,2],  
                    dtype=str)
```

String used to separate values  
Skip the first 2 lines  
Read the 1st and 3rd column  
The type of the resulting array

##### Files with mixed data types

```
>>> filename = 'titanic.csv'  
>>> data = np.genfromtxt(filename,  
                    delimiter=',',  
                    names=True,  
                    dtype=None)
```

Look for column header

```
>>> data_array = np.recfromcsv(filename)
```

The default `dtype` of the `np.recfromcsv()` function is `None`.

#### Importing Flat Files with pandas

```
>>> filename = 'winequality-red.csv'  
>>> data = pd.read_csv(filename,  
                    nrows=5,  
                    header=None,  
                    sep='\t',  
                    comment='#',  
                    na_values=[''])
```

Number of rows of file to read  
Row number to use as col names  
Delimiter to use  
Character to split comments  
String to recognize as NA/NaN

### Excel Spreadsheets

```
>>> file = 'urbanpop.xlsx'  
>>> data = pd.ExcelFile(file)  
>>> df_sheet2 = data.parse('1960-1966',  
                           skiprows=[0],  
                           names=['Country',  
                                  'AAM: War(2002)'])  
  
>>> df_sheet1 = data.parse(0,  
                           parse_cols=[0],  
                           skiprows=[0],  
                           names=['Country'])
```

To access the sheet names, use the `sheet_names` attribute:

```
>>> data.sheet_names
```

### SAS Files

```
>>> from sas7bdat import SAS7BDAT  
>>> with SAS7BDAT('urbanpop.sas/bdat') as file:  
    df_sas = file.to_data_frame()
```

### Stata Files

```
>>> data = pd.read_stata('urbanpop.dta')
```

### Relational Databases

```
>>> from sqlalchemy import create_engine  
>>> engine = create_engine('sqlite:///Northwind.sqlite')
```

Use the `table_names()` method to fetch a list of table names:

```
>>> table_names = engine.table_names()
```

#### Querying Relational Databases

```
>>> con = engine.connect()  
>>> rs = con.execute("SELECT * FROM Orders")  
>>> df = pd.DataFrame(rs.fetchall())  
>>> df.columns = rs.keys()  
>>> con.close()
```

#### Using the context manager with

```
>>> with engine.connect() as con:  
    rs = con.execute("SELECT OrderID FROM Orders")  
    df = pd.DataFrame(rs.fetchmany(size=5))  
    df.columns = rs.keys()
```

#### Querying relational databases with pandas

```
>>> df = pd.read_sql_query("SELECT * FROM Orders", engine)
```

### Exploring Your Data

#### NumPy Arrays

```
>>> data_array.dtype  
>>> data_array.shape  
>>> len(data_array)
```

Data type of array elements  
Array dimensions  
Length of array

#### pandas DataFrames

```
>>> df.head()  
>>> df.tail()  
>>> df.index  
>>> df.columns  
>>> df.info()  
>>> data_array = data.values
```

Return first DataFrame rows  
Return last DataFrame rows  
Describe index  
Describe DataFrame columns  
Info on DataFrame  
Convert a DataFrame to an a NumPy array

### Pickled Files

```
>>> import pickle  
>>> with open('pickled_fruit.pkl', 'rb') as file:  
    pickled_data = pickle.load(file)
```

### HDF5 Files

```
>>> import h5py  
>>> filename = 'H-H1_LOSC_4_v1-815411200-4096.hdf5'  
>>> data = h5py.File(filename, 'r')
```

### Matlab Files

```
>>> import scipy.io  
>>> filename = 'workspace.mat'  
>>> mat = scipy.io.loadmat(filename)
```

### Exploring Dictionaries

#### Accessing Elements with Functions

<pre>&gt;&gt;&gt; print(mat.keys()) &gt;&gt;&gt; for key in mat.keys():     print(key)</pre> <p>meta quality strain</p>	<p>Print dictionary keys Print dictionary keys</p>
<pre>&gt;&gt;&gt; pickled_data.values() &gt;&gt;&gt; print(mat.items())</pre> <p>Return dictionary values Returns items in list format of (key, value) tuple pairs</p>	

#### Accessing Data Items with Keys

<pre>&gt;&gt;&gt; for key in data['meta'].keys():     print(key)</pre> <p>Description DescriptionURL Detector Duration GRSstart Observatory Type UTCstart</p>	<p>Explore the HDF5 structure</p>
<pre>&gt;&gt;&gt; print(data['meta']['Description'].value)</pre> <p>Retrieve the value for a key</p>	

### Navigating Your FileSystem

#### Magic Commands

<pre>!ls %cd .. %pwd</pre>	<p>List directory contents of files and directories Change current working directory Return the current working directory path</p>
------------------------------------	--

#### os Library

<pre>&gt;&gt;&gt; import os &gt;&gt;&gt; path = "/usr/tmp" &gt;&gt;&gt; wd = os.getcwd() &gt;&gt;&gt; os.listdir(wd) &gt;&gt;&gt; os.chdir(path) &gt;&gt;&gt; os.rename("test1.txt", "test2.txt") &gt;&gt;&gt; os.remove("test1.txt") &gt;&gt;&gt; os.mkdir("newdir")</pre>	<p>Store the name of current directory in a string Output contents of the directory in a list Change current working directory Rename a file Delete an existing file Create a new directory</p>
---	---



# Python & Pylab Cheat Sheet

## Running

```
python3           standard python shell.  
ipython3          improved interactive shell.  
ipython3 --pylab  ipython including pylab  
python3 file.py   run file.py  
python3 -i file.py run file.py, stay in interactive mode
```

To quit use `exit()` or `[ctrl]+[d]`

## Getting Help

```
help()            interactive Help  
help(object)     help for object  
object?          ipython: help for object  
object??         ipython: extended help for object  
%magic           ipython: help on magic commands
```

## Import Syntax, e.g. for $\pi$

```
import numpy      use: numpy.pi  
import numpy as np use: np.pi  
from numpy import pi    use: pi  
from numpy import *    use: pi (use sparingly)
```

## Types

i = 1	Integer
f = 1.	Float
c = 1+2j	Complex
True/False	Boolean
'abc'	String
"abc"	String

with this:

c.real	1.0
c.imag	2.0
c.conjugate()	1-2j

## Operators

mathematics	
+	addition
-	subtraction
*	multiplication
/i	int division
/f	float division
**	power
%	modulo

comparison	
=	assign
==	equal
!=	unequal
<	less
<=	less-equal
>=	greater-equal
>	greater

## Basic Syntax

```
raw_input('foo')      read string from command-line  
class Foo(Object): ... class definition  
def bar(args): ... function/method definition  
if c: ... elif c: ... else: branching  
try: ... except Error: ... exception handling  
while cond: ... while loop  
for item in list: ... for loop  
[item for item in list] for loop, list notation
```

## Useful tools

```
pylint file.py       static code checker  
pydoc file          parse docstring to man-page  
python3 -m doctest  run examples in docstring  
file.py             run in debugger  
python3 -m pdb file.py
```

# NumPy & Friends

The following import statement is assumed:  
`from pylab import *`

## General Math

f: float, c: complex:	
abs(c)	absolute value of f or c
sign(c)	get sign of f or c
fix(f)	round towards 0
floor(f)	round towards -inf
ceil(f)	round towards +inf
f.round(p)	round f to p places
angle(c)	angle of complex number
sin(c)	sinus of argument
arcsin(c)	arcsin of argument
cos, tan,...	analogous

## Defining Lists, Arrays, Matrices

l: list, a: array:	basic list
[[1,2],[3,4,5]]	array from "rectangular" list
matrix([[1,2],[3,4]])	matrix from 2d-list
range(min, max, step)	integers in [min, max)
list(range(...))	list from range()
arange(min, max, step)	integer array in [min, max)
frange(min, max, step)	float array in [min, max]
linspace(min, max, num)	num samples in [min, max]
meshgrid(x,y)	create coord-matrices
zeros, ones, eye	generate special arrays

## Element Access

l[row][col]	list: basic access
l[min:max]	list: range access [min,max)
a[row,col] or a[row][col]	array: basic access
a[min:max,min:max]	array: range access [min,max)
a[list]	array: select indices in list
a[np.where(cond)]	array: select where cond true

## List/Array Properties

len(l)	size of first dim
a.size	total number of entries
a.ndim	number of dimensions
a.shape	size along dimensions
ravel(l) or a.ravel()	convert to 1-dim
a.flat	iterate all entries

## Matrix Operations

a: array, M: matrix:	element-wise product
a*a	dot product
dot(a,a) or M*M	cross product
cross(a,a)	inverted matrix
inv(a) or M.I	transposed matrix
transpose(a) or M.T	calculate determinante
det(a)	

## Statistics

sum(l,d) or a.sum(d)	sum elements along d
mean(l,d) or a.mean(d)	mean along d
std(l,d) or a.std(d)	standard deviation along d
min(l,d) or a.min(d)	minima along d
max(l,d) or a.max(d)	maxima along d

## Misc functions

loadtxt(file)	read values from file
polyval(coeff,xvals)	evaluate polynomial at xvals
roots(coeff)	find roots of polynomial
map(func,list)	apply func on each element of list

## Plotting

### Plot Types

plot(xvals, yvals, 'g+')	mark 3 points with green +
errorbar()	like plot with error bars
semilogx(), semilogx()	like plot, semi-log axis
loglog()	double logarithmic plot
polar(phi_vals, rvals)	plot in polar coordinates
hist(vals, n_bins)	create histogram from values
bar(low_edge, vals, width)	create bar-plot
contour(xvals,yvals,zvals)	create contour-plot

### Pylab Plotting Equivalences

figure()	fig = figure()
	ax = axes()
subplot(2,1,1)	ax = fig.add_subplot(2,1,1)
plot()	ax.plot()
errorbar()	ax.errorbar()
semilogx, ...	analogous
polar()	axes(polar=True) and ax.plot()
axis()	ax.set_xlim(), ax.set_ylim()
grid()	ax.grid()
title()	ax.set_title()
xlabel()	ax.set_xlabel()
legend()	ax.legend()
colorbar()	fig.colorbar(plot)

## Plotting 3D

```
from mpl_toolkits.mplot3d import Axes3D  
ax = fig.add_subplot(...,projection='3d')  
or ax = Axes3D(fig)          create 3d-axes object  
ax.plot(xvals, yvals, zvals) normal plot in 3d  
ax.plot_wireframe          wire mesh  
ax.plot_surface            colored surface
```

License: CC-by-sa  
Copyright: January 15, 2016, Nicola Chiapolini  
<http://www.physik.uzh.ch/~nchiapol>

# Python For Data Science Cheat Sheet

## NumPy Basics

Learn Python for Data Science Interactively at [www.DataCamp.com](http://www.DataCamp.com)



### NumPy

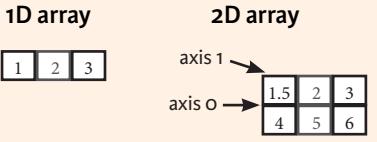
The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.

Use the following import convention:

```
>>> import numpy as np
```



### NumPy Arrays



### Creating Arrays

```
>>> a = np.array([1,2,3])
>>> b = np.array([(1.5,2,3), (4,5,6)], dtype = float)
>>> c = np.array([(1.5,2,3), (4,5,6)], [(3,2,1), (4,5,6)]),
      dtype = float)
```

### Initial Placeholders

```
>>> np.zeros((3,4))
>>> np.ones((2,3,4),dtype=np.int16)
>>> d = np.arange(10,25,5)

>>> np.linspace(0,2,9)

>>> e = np.full((2,2),7)
>>> f = np.eye(2)
>>> np.random.random((2,2))
>>> np.empty((3,2))
```

Create an array of zeros  
Create an array of ones  
Create an array of evenly spaced values (step value)  
Create an array of evenly spaced values (number of samples)  
Create a constant array  
Create a 2x2 identity matrix  
Create an array with random values  
Create an empty array

### I/O

#### Saving & Loading On Disk

```
>>> np.save('my_array', a)
>>> np.savetxt('array.npz', a, b)
>>> np.load('my_array.npy')
```

#### Saving & Loading Text Files

```
>>> np.loadtxt("myfile.txt")
>>> np.genfromtxt("my_file.csv", delimiter=',')
>>> np.savetxt("myarray.txt", a, delimiter=" ")
```

### Data Types

<code>&gt;&gt;&gt; np.int64</code>	Signed 64-bit integer types
<code>&gt;&gt;&gt; np.float32</code>	Standard double-precision floating point
<code>&gt;&gt;&gt; np.complex</code>	Complex numbers represented by 128 floats
<code>&gt;&gt;&gt; np.bool</code>	Boolean type storing TRUE and FALSE values
<code>&gt;&gt;&gt; np.object</code>	Python object type
<code>&gt;&gt;&gt; np.string_</code>	Fixed-length string type
<code>&gt;&gt;&gt; np_unicode_</code>	Fixed-length unicode type

### Inspecting Your Array

```
>>> a.shape
>>> len(a)
>>> b.ndim
>>> e.size
>>> b.dtype
>>> b.dtype.name
>>> b.astype(int)
```

Array dimensions  
Length of array  
Number of array dimensions  
Number of array elements  
Data type of array elements  
Name of data type  
Convert an array to a different type

### Asking For Help

```
>>> np.info(np.ndarray.dtype)
```

### Array Mathematics

#### Arithmetic Operations

```
>>> g = a - b
      array([[-0.5,  0. ,  0. ],
             [-3. , -3. , -3. ]])
>>> np.subtract(a,b)
>>> b + a
      array([[ 2.5,  4. ,  6. ],
             [ 5. ,  7. ,  9. ]])
>>> np.add(b,a)
>>> a / b
      array([[ 0.66666667,  1.        ,  1.        ],
             [ 0.25,  0.4,  0.5       ]])
>>> np.divide(a,b)
>>> a * b
      array([[ 1.5,  4. ,  9. ],
             [ 4. , 10. , 18. ]])
>>> np.multiply(a,b)
>>> np.exp(b)
>>> np.sqrt(b)
>>> np.sin(a)
>>> np.cos(b)
>>> np.log(a)
>>> e.dot(f)
      array([[ 7.,  7.],
             [ 7.,  7.]])
```

Subtraction  
Addition  
Addition  
Division  
Division  
Multiplication  
Multiplication  
Exponentiation  
Square root  
Print sines of an array  
Element-wise cosine  
Element-wise natural logarithm  
Dot product

### Comparison

```
>>> a == b
      array([[False,  True,  True],
             [False, False, False]], dtype=bool)
>>> a < 2
      array([True, False, False], dtype=bool)
>>> np.array_equal(a, b)
```

Element-wise comparison  
Element-wise comparison  
Array-wise comparison

### Aggregate Functions

```
>>> a.sum()
>>> a.min()
>>> b.max(axis=0)
>>> b.cumsum(axis=1)
>>> a.mean()
>>> b.median()
>>> a.correlcoef()
>>> np.std(b)
```

Array-wise sum  
Array-wise minimum value  
Maximum value of an array row  
Cumulative sum of the elements  
Mean  
Median  
Correlation coefficient  
Standard deviation

### Copying Arrays

```
>>> h = a.view()
>>> np.copy(a)
>>> h = a.copy()
```

Create a view of the array with the same data  
Create a copy of the array  
Create a deep copy of the array

### Sorting Arrays

```
>>> a.sort()
>>> c.sort(axis=0)
```

Sort an array  
Sort the elements of an array's axis

### Subsetting, Slicing, Indexing

#### Subsetting

```
>>> a[2]
      3
>>> b[1,2]
      6.0
```

Select the element at the 2nd index

#### Slicing

```
>>> a[0:2]
      array([1, 2])
>>> b[0:2,1]
      array([ 2.,  5.])
```

Select items at index 0 and 1

```
>>> b[:1]
      array([[1.5, 2., 3.]])
```

Select all items at row 0  
(equivalent to `b[0:1, :]`)  
Same as `[1, :, :]`

```
>>> c[1,:]
      array([[ 3.,  2.,  1.],
             [ 4.,  5.,  6.]])
```

Reversed array `a`

#### Boolean Indexing

```
>>> a[a<2]
      array([1])
```

Select elements `(1,0),(0,1),(1,2)` and `(0,0)`

#### Fancy Indexing

```
>>> b[[1, 0, 1, 0], [0, 1, 2, 0]]
      array([ 4.,  2.,  6., 1.5])
>>> b[[1, 0, 1, 0]][:, [0, 1, 2, 0]]
      array([[ 4.,  5.,  6.,  4.],
             [ 1.5,  2.,  3.,  1.5],
             [ 4.,  5.,  6.,  4.],
             [ 1.5,  2.,  3.,  1.5]])
```

Select a subset of the matrix's rows and columns

### Array Manipulation

#### Transposing Array

```
>>> i = np.transpose(b)
>>> i.T
```

Permute array dimensions  
Permute array dimensions

#### Changing Array Shape

```
>>> b.ravel()
>>> g.reshape(3,-2)
```

Flatten the array  
Reshape, but don't change data

#### Adding/Removing Elements

```
>>> h.resize((2,6))
>>> np.append(h,g)
>>> np.insert(a, 1, 5)
>>> np.delete(a,[1])
```

Return a new array with shape `(2,6)`  
Append items to an array  
Insert items in an array  
Delete items from an array

#### Combining Arrays

```
>>> np.concatenate((a,d),axis=0)
      array([ 1,  2,  3, 10, 15, 20])
>>> np.vstack((a,b))
      array([[ 1.,  2.,  3.],
             [ 1.5,  2.,  3.],
             [ 4.,  5.,  6.]])
>>> np.r_[e,f]
>>> np.hstack((e,f))
      array([[ 7.,  7.,  1.,  0.],
             [ 7.,  7.,  0.,  1.]])
>>> np.column_stack((a,d))
      array([[ 1, 10],
             [ 2, 15],
             [ 3, 20]])
>>> np.c_[a,d]
```

Concatenate arrays  
Stack arrays vertically (row-wise)  
Stack arrays vertically (row-wise)  
Stack arrays horizontally (column-wise)

Create stacked column-wise arrays  
Create stacked column-wise arrays

Create stacked column-wise arrays  
Split the array horizontally at the 3rd index  
Split the array vertically at the 2nd index



# Python For Data Science Cheat Sheet

## SciPy - Linear Algebra

Learn More Python for Data Science [Interactively](#) at [www.datacamp.com](http://www.datacamp.com)



### SciPy

The SciPy library is one of the core packages for scientific computing that provides mathematical algorithms and convenience functions built on the NumPy extension of Python.



### Interacting With NumPy

[Also see NumPy](#)

```
>>> import numpy as np
>>> a = np.array([1,2,3])
>>> b = np.array([(1+5j,2j,3j), (4j,5j,6j)])
>>> c = np.array([(1.5,2,3), (4,5,6)], [(3,2,1), (4,5,6)]])
```

#### Index Tricks

>>> np.mgrid[0:5,0:5]	Create a dense meshgrid
>>> np.ogrid[0:2,0:2]	Create an open meshgrid
>>> np.r_[3,[0]*5,-1:1:10j]	Stack arrays vertically (row-wise)
>>> np.c_[[b,c]]	Create stacked column-wise arrays

#### Shape Manipulation

>>> np.transpose(b)	Permute array dimensions
>>> b.flatten()	Flatten the array
>>> np.hstack((b,c))	Stack arrays horizontally (column-wise)
>>> np.vstack((a,b))	Stack arrays vertically (row-wise)
>>> np.hsplit(c,2)	Split the array horizontally at the 2nd index
>>> np.vsplit(d,2)	Split the array vertically at the 2nd index

#### Polynomials

>>> from numpy import poly1d	
>>> p = poly1d([3,4,5])	Create a polynomial object

#### Vectorizing Functions

>>> def myfunc(a): if a < 0: return a**2 else: return a/2	
>>> np.vectorize(myfunc)	Vectorize functions

#### Type Handling

>>> np.real(b)	Return the real part of the array elements
>>> np.imag(b)	Return the imaginary part of the array elements
>>> np.real_if_close(c,tol=1000)	Return a real array if complex parts close to 0
>>> np.cast['f'](np.pi)	Cast object to a data type

#### Other Useful Functions

>>> np.angle(b,deg=True)	Return the angle of the complex argument
>>> g = np.linspace(0,np.pi,num=5)	Create an array of evenly spaced values (number of samples)
>>> g[3:] += np.pi	Unwrap
>>> np.unwrap(g)	Create an array of evenly spaced values (log scale)
>>> np.logspace(0,10,3)	Return values from a list of arrays depending on conditions
>>> np.select([c<4],[c*2])	Factorial
>>> misc.factorial(a)	Combine N things taken at k time
>>> misc.comb(10,3,exact=True)	Weights for N-point central derivative
>>> misc.central_diff_weights(3)	Find the n-th derivative of a function at a point
>>> misc.derivative(myfunc,1.0)	

## Linear Algebra

You'll use the `linalg` and `sparse` modules. Note that `scipy.linalg` contains and expands on `numpy.linalg`.

```
>>> from scipy import linalg, sparse
```

### Creating Matrices

```
>>> A = np.matrix(np.random.random((2,2)))
>>> B = np.asmatrix(b)
>>> C = np.mat(np.random.random((10,5)))
>>> D = np.mat([(3,4), (5,6)])
```

### Basic Matrix Routines

#### Inverse

```
>>> A.I
```

#### Transposition

```
>>> A.T
```

#### Trace

```
>>> np.trace(A)
```

#### Norm

```
>>> linalg.norm(A)
>>> linalg.norm(A,1)
```

```
>>> linalg.norm(A,np.inf)
```

#### Rank

```
>>> np.linalg.matrix_rank(C)
```

#### Determinant

```
>>> linalg.det(A)
```

#### Solving linear problems

```
>>> linalg.solve(A,b)
```

```
>>> E = np.mat(a).T
```

```
>>> linalg.lstsq(F,E)
```

#### Generalized inverse

```
>>> linalg.pinv(C)
```

```
>>> linalg.pinv2(C)
```

#### Inverse

#### Inverse

#### Transpose matrix

#### Conjugate transposition

#### Trace

#### Frobenius norm

L<sub>1</sub> norm (max column sum)

L<sub>inf</sub> norm (max row sum)

#### Matrix rank

#### Determinant

Solver for dense matrices  
Solver for dense matrices  
Least-squares solution to linear matrix equation

Compute the pseudo-inverse of a matrix  
(least-squares solver)

Compute the pseudo-inverse of a matrix  
(SVD)

### Creating Sparse Matrices

```
>>> F = np.eye(3, k=1)
>>> G = np.mat(np.identity(2))
```

Create a 2x2 identity matrix

Create a 2x2 identity matrix

```
>>> C[> 0.5] = 0
```

Compressed Sparse Row matrix

Compressed Sparse Column matrix

Dictionary Of Keys matrix

Sparse matrix to full matrix

Identify sparse matrix

### Sparse Matrix Routines

#### Inverse

#### Inverse

#### Norm

#### Norm

#### Solving linear problems

Solver for sparse matrices

### Sparse Matrix Functions

```
>>> sparse.linalg.expm(I)
```

Sparse matrix exponential

### Asking For Help

```
>>> help(scipy.linalg.diagsvd)
>>> np.info(np.matrix)
```

[Also see NumPy](#)

### Matrix Functions

#### Addition

```
>>> np.add(A,D)
```

#### Subtraction

```
>>> np.subtract(A,D)
```

#### Division

```
>>> np.divide(A,D)
```

#### Multiplication

```
>>> A @ D
```

#### Exponential Functions

```
>>> np.multiply(D,A)
```

```
>>> np.dot(A,D)
```

```
>>> np.vdot(A,D)
```

```
>>> np.inner(A,D)
```

```
>>> np.outer(A,D)
```

```
>>> np.tensordot(A,D)
```

```
>>> np.kron(A,D)
```

#### Logarithm Function

```
>>> np.logm(A)
```

#### Trigonometric Functions

```
>>> np.sinm(D)
```

```
>>> np.cosm(D)
```

```
>>> np.tanm(A)
```

#### Hyperbolic Trigonometric Functions

```
>>> np.sinhm(D)
```

```
>>> np.coshm(D)
```

```
>>> np.tanhm(A)
```

#### Matrix Sign Function

```
>>> np.signm(A)
```

#### Matrix Square Root

```
>>> np.sqrtm(A)
```

#### Arbitrary Functions

```
>>> linalg.funm(A, lambda x: x*x)
```

#### Addition

#### Subtraction

#### Division

Multiplication operator  
(Python 3)  
Multiplication  
Dot product  
Vector dot product

Inner product  
Outer product  
Tensor dot product  
Kronecker product

Matrix exponential  
Matrix exponential (Taylor Series)  
Matrix exponential (eigenvalue decomposition)

#### Matrix logarithm

Matrix sine  
Matrix cosine  
Matrix tangent

Hyperbolic matrix sine  
Hyperbolic matrix cosine  
Hyperbolic matrix tangent

#### Matrix sign function

#### Matrix square root

#### Evaluate matrix function

### Decompositions

#### Eigenvalues and Eigenvectors

```
>>> la, v = linalg.eig(A)
```

Solve ordinary or generalized eigenvalue problem for square matrix  
Unpack eigenvalues

First eigenvector

Second eigenvector

Unpack eigenvalues

Singular Value Decomposition (SVD)

Construct sigma matrix in SVD

#### LU Decomposition

```
>>> P,L,U = linalg.lu(C)
```

LU Decomposition

### Sparse Matrix Decompositions

#### Eigenvalues and Eigenvectors SVD

```
>>> la, v = sparse.linalg.eigs(F, 1)
```

```
>>> sparse.linalg.svds(H, 2)
```

Eigenvalues and eigenvectors SVD

DataCamp

Learn Python for Data Science [Interactively](#)



# Python For Data Science Cheat Sheet

## Pandas Basics

Learn Python for Data Science Interactively at [www.DataCamp.com](http://www.DataCamp.com)



### Pandas

The Pandas library is built on NumPy and provides easy-to-use data structures and data analysis tools for the Python programming language.



Use the following import convention:

```
>>> import pandas as pd
```

### Pandas Data Structures

#### Series

A one-dimensional labeled array capable of holding any data type

a	3
b	-5
c	7
d	4

Index →

```
>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```

#### DataFrame

Index	Columns		
	Country	Capital	Population
0	Belgium	Brussels	11190846
1	India	New Delhi	1303171035
2	Brazil	Brasilia	207847528

A two-dimensional labeled data structure with columns of potentially different types

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],
   >>>          'Capital': ['Brussels', 'New Delhi', 'Brasilia'],
   >>>          'Population': [11190846, 1303171035, 207847528]}
>>> df = pd.DataFrame(data,
   >>>                      columns=['Country', 'Capital', 'Population'])
```

### I/O

#### Read and Write to CSV

```
>>> pd.read_csv('file.csv', header=None, nrows=5)
>>> df.to_csv('myDataFrame.csv')
```

#### Read and Write to Excel

```
>>> pd.read_excel('file.xlsx')
>>> pd.to_excel('dir/myDataFrame.xlsx', sheet_name='Sheet1')
Read multiple sheets from the same file
>>> xlsx = pd.ExcelFile('file.xls')
>>> df = pd.read_excel(xlsx, 'Sheet1')
```

### Asking For Help

```
>>> help(pd.Series.loc)
```

### Selection

#### Getting

```
>>> s['b']
-5
>>> df[1:]
   Country    Capital  Population
1  India      New Delhi     1303171035
2  Brazil     Brasilia     207847528
```

Also see NumPy Arrays

Get one element

Get subset of a DataFrame

### Selecting, Boolean Indexing & Setting

#### By Position

```
>>> df.iloc[[0], [0]]
'Belgium'
>>> df.iat[[0], [0]]
'Belgium'
```

#### By Label

```
>>> df.loc[[0], ['Country']]
'Belgium'
>>> df.at[[0], ['Country']]
'Belgium'
```

#### By Label/Position

```
>>> df.ix[2]
   Country      Brazil
   Capital    Brasilia
   Population  207847528
```

```
>>> df.ix[:, 'Capital']
0    Brussels
1   New Delhi
2    Brasilia
```

```
>>> df.ix[1, 'Capital']
'New Delhi'
```

#### Boolean Indexing

```
>>> s[~(s > 1)]
>>> s[(s < -1) | (s > 2)]
>>> df[df['Population'] > 12000000000]
```

#### Setting

```
>>> s['a'] = 6
```

Select single value by row & column

Select single value by row & column labels

Select single row of subset of rows

Select a single column of subset of columns

Select rows and columns

Series s where value is not >1  
s where value is <-1 or >2  
Use filter to adjust DataFrame

Set index a of Series s to 6

### Dropping

```
>>> s.drop(['a', 'c'])
>>> df.drop('Country', axis=1)
```

Drop values from rows (axis=0)

Drop values from columns (axis=1)

### Sort & Rank

```
>>> df.sort_index()
>>> df.sort_values(by='Country')
>>> df.rank()
```

Sort by labels along an axis  
Sort by the values along an axis  
Assign ranks to entries

### Retrieving Series/DataFrame Information

#### Basic Information

```
>>> df.shape
>>> df.index
>>> df.columns
>>> df.info()
>>> df.count()
```

(rows,columns)  
Describe index  
Describe DataFrame columns  
Info on DataFrame  
Number of non-NA values

#### Summary

```
>>> df.sum()
>>> df.cumsum()
>>> df.min() / df.max()
>>> df.idxmin() / df.idxmax()
>>> df.describe()
>>> df.mean()
>>> df.median()
```

Sum of values  
Cummulative sum of values  
Minimum/maximum values  
Minimum/Maximum index value  
Summary statistics  
Mean of values  
Median of values

### Applying Functions

```
>>> f = lambda x: x**2
>>> df.apply(f)
>>> df.applymap(f)
```

Apply function  
Apply function element-wise

### Data Alignment

#### Internal Data Alignment

NA values are introduced in the indices that don't overlap:

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
>>> s + s3
a    10.0
b    NaN
c     5.0
d     7.0
```

### Arithmetic Operations with Fill Methods

You can also do the internal data alignment yourself with the help of the fill methods:

```
>>> s.add(s3, fill_value=0)
a    10.0
b    -5.0
c     5.0
d     7.0
>>> s.sub(s3, fill_value=2)
>>> s.div(s3, fill_value=4)
>>> s.mul(s3, fill_value=3)
```



## Summarize Data

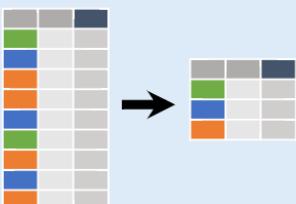
```
df['w'].value_counts()
Count number of rows with each unique value of variable
len(df)
# of rows in DataFrame.
df['w'].nunique()
# of distinct values in a column.
df.describe()
Basic descriptive statistics for each column (or GroupBy)
```



pandas provides a large set of **summary functions** that operate on different kinds of pandas objects (DataFrame columns, Series, GroupBy, Expanding and Rolling (see below)) and produce single values for each of the groups. When applied to a DataFrame, the result is returned as a pandas Series for each column. Examples:

<b>sum()</b>	<b>min()</b>
Sum values of each object.	Minimum value in each object.
<b>count()</b>	<b>max()</b>
Count non-NA/null values of each object.	Maximum value in each object.
<b>median()</b>	<b>mean()</b>
Median value of each object.	Mean value of each object.
<b>quantile([0.25, 0.75])</b>	<b>var()</b>
Quantiles of each object.	Variance of each object.
<b>std()</b>	<b>std()</b>
Apply function to each object.	Standard deviation of each object.

## Group Data



**df.groupby(by="col")**  
Return a GroupBy object, grouped by values in column named "col".

**df.groupby(level="ind")**  
Return a GroupBy object, grouped by values in index level named "ind".

All of the summary functions listed above can be applied to a group.

Additional GroupBy functions:

<b>size()</b>	<b>agg(function)</b>
Size of each group.	Aggregate group using function.

## Windows

```
df.expanding()
Return an Expanding object allowing summary functions to be applied cumulatively.
df.rolling(n)
Return a Rolling object allowing summary functions to be applied to windows of length n.
```

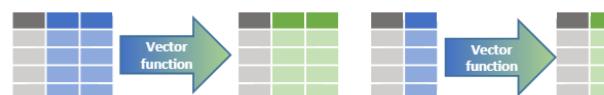
## Handling Missing Data

```
df.dropna()
Drop rows with any column having NA/null data.
df.fillna(value)
Replace all NA/null data with value.
```

## Make New Columns



```
df.assign(Area=lambda df: df.Length*df.Height)
Compute and append one or more new columns.
df['Volume'] = df.Length*df.Height*df.Depth
Add single column.
pd.qcut(df.col, n, labels=False)
Bin column into n buckets.
```



pandas provides a large set of **vector functions** that operate on all columns of a DataFrame or a single selected column (a pandas Series). These functions produce vectors of values for each of the columns, or a single Series for the individual Series. Examples:

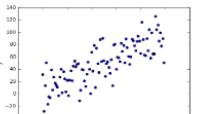
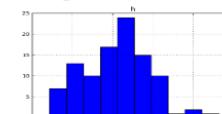
<b>max(axis=1)</b>	<b>min(axis=1)</b>
Element-wise max.	Element-wise min.
<b>clip(lower=-10,upper=10)</b>	<b>abs()</b>
Trim values at input thresholds	Absolute value.

The examples below can also be applied to groups. In this case, the function is applied on a per-group basis, and the returned vectors are of the length of the original DataFrame.

<b>shift(1)</b>	<b>shift(-1)</b>
Copy with values shifted by 1.	Copy with values lagged by 1.
<b>rank(method='dense')</b>	<b>cumsum()</b>
Ranks with no gaps.	Cumulative sum.
<b>rank(method='min')</b>	<b>cummax()</b>
Ranks. Ties get min rank.	Cumulative max.
<b>rank(pct=True)</b>	<b>cummin()</b>
Ranks rescaled to interval [0, 1].	Cumulative min.
<b>rank(method='first')</b>	<b>cumprod()</b>
Ranks. Ties go to first value.	Cumulative product.

## Plotting

```
df.plot.hist()
Histogram for each column
df.plot.scatter(x='w',y='h')
Scatter chart using pairs of points
```



## Combine Data Sets

adf		bdf	
x1	x2	x1	x3
A	1	A	T
B	2	B	F
C	3	D	T

+	=
---	---

### Standard Joins

x1	x2	x3
A	1	T
B	2	F
C	3	NaN

```
pd.merge(adf, bdf,
how='left', on='x1')
Join matching rows from bdf to adf.
```

x1	x2	x3
A	1.0	T
B	2.0	F
D	NaN	T

```
pd.merge(adf, bdf,
how='right', on='x1')
Join matching rows from adf to bdf.
```

x1	x2	x3
A	1	T
B	2	F

```
pd.merge(adf, bdf,
how='inner', on='x1')
Join data. Retain only rows in both sets.
```

x1	x2	x3
A	1	T
B	2	F
C	3	NaN
D	NaN	T

```
pd.merge(adf, bdf,
how='outer', on='x1')
Join data. Retain all values, all rows.
```

### Filtering Joins

x1	x2
A	1
B	2

```
adf[adf.x1.isin(bdf.x1)]
All rows in adf that have a match in bdf.
```

x1	x2
C	3

```
adf[~adf.x1.isin(bdf.x1)]
All rows in adf that do not have a match in bdf.
```

ydf		zdf	
x1	x2	x1	x2
A	1	B	2
B	2	C	3
C	3	D	4

### Set-like Operations

x1	x2
B	2
C	3

```
pd.merge(ydf, zdf)
Rows that appear in both ydf and zdf
(Intersection).
```

x1	x2
A	1
B	2
C	3
D	4

```
pd.merge(ydf, zdf, how='outer')
Rows that appear in either or both ydf and zdf
(Union).
```

x1	x2
A	1

```
pd.merge(ydf, zdf, how='outer',
indicator=True)
.query('_merge == "left_only"')
.drop(['_merge'], axis=1)
Rows that appear in ydf but not zdf (Setdiff).
```

# Python For Data Science Cheat Sheet

## Bokeh

Learn Bokeh [Interactively at www.DataCamp.com](#), taught by Bryan Van de Ven, core contributor

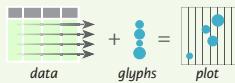


### Plotting With Bokeh

The Python interactive visualization library Bokeh enables high-performance visual presentation of large datasets in modern web browsers.



Bokeh's mid-level general purpose `bokeh.plotting` interface is centered around two main components: data and glyphs.



The basic steps to creating plots with the `bokeh.plotting` interface are:

1. Prepare some data:  
Python lists, NumPy arrays, Pandas DataFrames and other sequences of values
2. Create a new plot
3. Add renderers for your data, with visual customizations
4. Specify where to generate the output
5. Show or save the results

```
>>> from bokeh.plotting import figure
>>> from bokeh.io import output_file, show
>>> x = [1, 2, 3, 4, 5]          Step 1
>>> y = [6, 7, 2, 4, 5]
>>> p = figure(title="simple line example",
    x_axis_label="x",
    y_axis_label="y")
>>> p.line(x, y, legend="Temp.", line_width=2)  Step 2
>>> output_file("lines.html")   Step 3
>>> show(p)                  Step 4
>>> Step 5
```

## 1 Data

[Also see Lists, NumPy & Pandas](#)

Under the hood, your data is converted to Column Data Sources. You can also do this manually:

```
>>> import numpy as np
>>> import pandas as pd
>>> df = pd.DataFrame(np.array([[33, 9, 4, 65, 'US'],
    [32, 4, 4, 66, 'Asia'],
    [21, 4, 4, 109, 'Europe']]),
    columns=['mpg', 'cyl', 'hp', 'origin'],
    index=['Toyota', 'Fiat', 'Volvo'])
```

## 2 Plotting

```
>>> from bokeh.plotting import figure
>>> p1 = figure(plot_width=300, tools='pan,box_zoom')
>>> p2 = figure(plot_width=300, plot_height=300,
    x_range=(0, 8), y_range=(0, 8))
>>> p3 = figure()
```

## 3 Renderers & Visual Customizations

### Glyphs

**Scatter Markers**  
`>>> p1.circle(np.array([1,2,3]), np.array([3,2,1]),
 fill_color='white')`  
`>>> p2.square(np.array([1.5,3.5,5.5]), [1,4,3],
 color='blue', size=1)`

**Line Glyphs**  
`>>> p1.line([1,2,3,4], [3,4,5,6], line_width=2)`  
`>>> p2.multi_line(pd.DataFrame([[1,2,3],[5,6,7]]),
 pd.DataFrame([[3,4,5],[3,2,1]]), color="blue")`

### Rows & Columns Layout

**Rows**  
`>>> from bokeh.layouts import row`  
`>>> layout = row(p1,p2,p3)`

**Columns**  
`>>> from bokeh.layouts import column`  
`>>> layout = column(p1,p2,p3)`

**Nesting Rows & Columns**  
`>>> layout = row(column(p1,p2), p3)`

### Grid Layout

```
>>> from bokeh.layouts import gridplot
>>> row1 = [p1,p2]
>>> row2 = [p3]
>>> layout = gridplot([[p1,p2],[p3]])
```

### Tabbed Layout

```
>>> from bokeh.models.widgets import Panel, Tabs
>>> tab1 = Panel(child=p1, title="tab1")
>>> tab2 = Panel(child=p2, title="tab2")
>>> layout = Tabs(tabs=[tab1, tab2])
```

### Legends

#### Legend Location

**Inside Plot Area**  
`>>> p.legend.location = 'bottom_left'`

**Outside Plot Area**  
`>>> r1 = p2.asterisk(np.array([1,2,3]), np.array([3,2,1]))
 r2 = p2.line([1,2,3,4], [3,4,5,6])
 legend = Legend(items=[("One", [p1, r1]), ("Two", [r2])], location=(0, -30))
 p.add_layout(legend, 'right')`

## 4 Output

### Output to HTML File

```
>>> from bokeh.io import output_file, show
>>> output_file('my_bar_chart.html', mode='cdn')
```

### Notebook Output

```
>>> from bokeh.io import output_notebook, show
>>> output_notebook()
```

### Embedding

**Standalone HTML**  
`>>> from bokeh.embed import file_html`  
`>>> html = file_html(p, CDN, "my_plot")`

**Components**  
`>>> from bokeh.embed import components`  
`>>> script, div = components(p)`

## 5 Show or Save Your Plots

```
>>> show(p1)           >>> save(p1)
>>> show(layout)       >>> save(layout)
```

### Customized Glyphs

[Also see Data](#)  
**Selection and Non-Selection Glyphs**  
`>>> p = figure(tools='box_select')`  
`>>> p.circle('mpg', 'cyl', source=cds_df,
 selection_color='red',
 nonselection_alpha=0.1)`

**Hover Glyphs**  
`>>> hover = HoverTool(tooltips=None, mode='vline')`  
`>>> p3.add_tools(hover)`

**Colormapping**  
`>>> color_mapper = CategoricalColorMapper(
 factors=['US', 'Asia', 'Europe'],
 palette=['blue', 'red', 'green'])`  
`>>> p3.circle('mpg', 'cyl', source=cds_df,
 color=dict(field='origin',
 transform=color_mapper),
 legend='Origin')`

[Also see Data](#)

### Linked Plots

**Linked Axes**  
`>>> p2.x_range = p1.x_range`  
`>>> p2.y_range = p1.y_range`

**Linked Brushing**  
`>>> p4 = figure(plot_width = 100, tools='box_select,lasso_select')`  
`>>> p4.circle('mpg', 'cyl', source=cds_df)`  
`>>> p5 = figure(plot_width = 200, tools='box_select,lasso_select')`  
`>>> p5.circle('mpg', 'hp', source=cds_df)`  
`>>> layout = row(p4,p5)`

[Also see Data](#)

### Legend Orientation

```
>>> p.legend.orientation = "horizontal"
>>> p.legend.orientation = "vertical"
```

### Legend Background & Border

```
>>> p.legend.border_line_color = "navy"
>>> p.legend.background_fill_color = "white"
```

## Statistical Charts With Bokeh

[Also see Data](#)  
Bokeh's high-level `bokeh.charts` interface is ideal for quickly creating statistical charts

### Bar Chart

```
>>> from bokeh.charts import Bar
>>> p = Bar(df, stacked=True, palette=['red','blue'])
```

### Box Plot

```
>>> from bokeh.charts import BoxPlot
>>> p = BoxPlot(df, values='vals', label='cyl',
    legend='bottom_right')
```

### Histogram

```
>>> from bokeh.charts import Histogram
>>> p = Histogram(df, title='Histogram')
```

### Scatter Plot

```
>>> from bokeh.charts import Scatter
>>> p = Scatter(df, x='mpg', y='hp', marker='square',
    xlabel='Miles Per Gallon',
    ylabel='Horsepower')
```

**DataCamp**  
Learn Python for Data Science [Interactively](#)



# Python For Data Science Cheat Sheet

## Seaborn

Learn Data Science Interactively at [www.DataCamp.com](http://www.DataCamp.com)



### Statistical Data Visualization With Seaborn

The Python visualization library **Seaborn** is based on matplotlib and provides a high-level interface for drawing attractive statistical graphics.

Make use of the following aliases to import the libraries:

```
>>> import matplotlib.pyplot as plt  
>>> import seaborn as sns
```

The basic steps to creating plots with Seaborn are:

1. Prepare some data
2. Control figure aesthetics
3. Plot with Seaborn
4. Further customize your plot

```
>>> import matplotlib.pyplot as plt  
>>> import seaborn as sns  
>>> tips = sns.load_dataset("tips")  
>>> sns.set_style("whitegrid")  
>>> g = sns.lmplot(x="tip",  
y="total_bill",  
data=tips,  
aspects=2)  
>>> g = (g.set_axis_labels("Tip", "Total bill (USD)")  
set(xlim=(0,10), ylim=(0,100))  
>>> plt.title("title")  
>>> plt.show(g)
```

Step 1  
Step 2  
Step 3  
Step 4  
Step 5

## 1) Data

Also see Lists, NumPy & Pandas

```
>>> import pandas as pd  
>>> import numpy as np  
>>> uniform_data = np.random.rand(10, 12)  
>>> data = pd.DataFrame({'x':np.arange(1,101),  
'y':np.random.normal(0,4,100)})
```

Seaborn also offers built-in data sets:

```
>>> titanic = sns.load_dataset("titanic")  
>>> iris = sns.load_dataset("iris")
```

## 3) Plotting With Seaborn

### Axis Grids

```
>>> g = sns.FacetGrid(titanic,  
col="survived",  
row="sex")  
>>> g = g.map(plt.hist, "age")  
>>> sns.factorplot(x="pclass",  
y="survived",  
hue="sex",  
data=titanic)  
>>> sns.lmplot(x="sepal_width",  
y="sepal_length",  
hue="species",  
data=iris)
```

Subplot grid for plotting conditional relationships

Draw a categorical plot onto a Facetgrid

Plot data and regression model fits across a FacetGrid

```
>>> h = sns.PairGrid(iris)  
>>> h = h.map(plt.scatter)  
>>> sns.pairplot(iris)  
>>> i = sns.JointGrid(x="x",  
y="y",  
data=data)  
>>> i = i.plot(sns.regplot,  
sns.distplot)  
>>> sns.jointplot("sepal_length",  
"sepal_width",  
data=iris,  
kind="kde")
```

Subplot grid for plotting pairwise relationships  
Plot pairwise bivariate distributions  
Grid for bivariate plot with marginal univariate plots

Plot bivariate distribution

### Categorical Plots

```
>>> sns.stripplot(x="species",  
y="petal_length",  
data=iris)  
>>> sns.swarmplot(x="species",  
y="petal_length",  
data=iris)  
>>> sns.barplot(x="sex",  
y="survived",  
hue="class",  
data=titanic)  
>>> sns.countplot(x="deck",  
data=titanic,  
palette="Greens_d")  
>>> sns.pointplot(x="class",  
y="survived",  
hue="sex",  
data=titanic,  
palette=({"male": "g",  
"female": "m"},  
markers=["^", "o"],  
linestyles=[ "--", "-"])
```

Scatterplot with one categorical variable

Categorical scatterplot with non-overlapping points

Show point estimates and confidence intervals with scatterplot glyphs

Show count of observations

Show point estimates and confidence intervals as rectangular bars

```
>>> sns.boxplot(x="alive",  
y="age",  
hue="adult_male",  
data=titanic)  
>>> sns.boxplot(data=iris, orient="h")  
>>> sns.violinplot(x="age",  
y="sex",  
hue="survived",  
data=titanic)
```

Boxplot

Boxplot with wide-form data

Violin plot

### Regression Plots

```
>>> sns.regplot(x="sepal_width",  
y="sepal_length",  
data=iris,  
ax=ax)
```

Plot data and a linear regression model fit

### Distribution Plots

```
>>> plot = sns.distplot(data.y,  
kde=False,  
color="b")
```

Plot univariate distribution

### Matrix Plots

```
>>> sns.heatmap(uniform_data, vmin=0, vmax=1)
```

Heatmap

## 4) Further Customizations

Also see Matplotlib

### Axisgrid Objects

```
>>> g.despine(left=True)  
>>> g.set_ylabels("Survived")  
>>> g.set_xticklabels(rotation=45)  
>>> g.set_axis_labels("Survived",  
"Sex")  
>>> h.set(xlim=(0,5),  
ylim=(0,5),  
xticks=[0,2.5,5],  
yticks=[0,2.5,5])
```

Remove left spine  
Set the labels of the y-axis  
Set the tick labels for x  
Set the axis labels

Set the limit and ticks of the x-and y-axis

### Plot

```
>>> plt.title("A Title")  
>>> plt.ylabel("Survived")  
>>> plt.xlabel("Sex")  
>>> plt.ylim(0,100)  
>>> plt.xlim(0,10)  
>>> plt.set(ax, yticks=[0,5])  
>>> plt.tight_layout()
```

Add plot title  
Adjust the label of the y-axis  
Adjust the label of the x-axis  
Adjust the limits of the y-axis  
Adjust the limits of the x-axis  
Adjust a plot property  
Adjust subplot params

## 5) Show or Save Plot

Also see Matplotlib

```
>>> plt.show()  
>>> plt.savefig("foo.png")  
>>> plt.savefig("foo.png",  
transparent=True)
```

Show the plot  
Save the plot as a figure  
Save transparent figure

### Close & Clear

Also see Matplotlib

```
>>> plt.clf()  
>>> plt.cla()  
>>> plt.close()
```

Clear an axis  
Clear an entire figure  
Close a window





## GETTING STARTED

### 1. Install

In the terminal  
sudo pip install plotly

### 2. Sign Up & Configure

<http://www.plot.ly/python/getting-started>

### 3. Boilerplate Imports

```
import plotly.plotly as py
import plotly.graph_objs as go
```

### 4. A Hello World Figure

```
trace = {'x': [1, 2], 'y': [1, 2]}
data = [trace]
fig = go.Figure()
    data = data, layout = layout )
```

### 5. Plot the Figure!

In the terminal:  
plot\_url = py.plot ( fig )

Or in the IPython notebook:  
py.iplot ( fig )

## BASIC CHARTS

### Line Plots

```
trace1 = go.Scatter (
    x = [ 1, 2 ], y = [ 1, 2 ])
trace2 = go.Scatter (
    x = [ 1, 2 ], y = [ 2, 1 ])
py.iplot ([ trace1, trace2 ])
```

### Bubble Charts

```
trace = go.Scatter (
    x = [ 1, 2, 3 ], y = [ 1, 2, 3 ],
    marker = dict (
        color = [ 'red', 'blue',
        'green' ],
        size = [ 30, 80, 200 ],
        mode = 'markers' )
    py.iplot ([ trace ])
```

### Scatter Plots

```
trace1 = go.Scatter (
    x = [ 1, 2, 3 ], y = [ 1, 2, 3 ],
    text = [ 'A', 'B', 'C' ],
    textposition = 'top center',
    mode = 'markers+text' )
mode = [ trace ]
py.iplot ( data )
```

### Heatmaps

```
trace = go.Heatmap (
    z = [ [ 1, 2, 3, 4 ],
        [ 5, 6, 7, 8 ] ],
    data = [ trace ]
    py.iplot ( data )
```

### Bar Charts

```
trace = go.Bar (
    x = [ 1, 2 ], y = [ 1, 2 ])
data = [ trace ]
py.iplot ( data )
```

### Area Plots

```
trace = go.Scatter (
    x = [ 1, 2 ], y = [ 1, 2 ],
    fill = 'tonexty' )
data = [ trace ]
py.iplot ( data )
```

## LAYOUT

### Legends

```
trace1 = go.Scatter (
    name = 'Calvin'
    x = [ 1, 2 ], y = [ 1, 2 ])
```

```
trace2 = go.Scatter (
    name = 'Hobbes'
    x = [ 2, 1 ], y = [ 2, 1 ])
```

```
layout = go.Layout (
    showlegend = True ,
    legend = dict (
        x = 0.2 , y = 0.5 )
    )
```

```
data = [ trace1 , trace2 ]
fig = go.Figure (
    data = data ,
    layout = layout )
py.iplot ( fig )
```

### Axes

```
trace = go.Scatter (
    x = [ 1, 2, 3, 4 ],
    y = [ 1, 2, 3, 6 ])
```

```
axis_template = dict (
    showgrid = False ,
    zeroline = False ,
    nticks = 20 ,
    showline = True ,
    title = 'X AXIS'
    mirror = 'all' )
layout = go.Layout (
    xaxis = axis_template ,
    yaxis = axis_template ,
    )
```

```
data = [ trace ]
fig = go.Figure (
    data = data
    layout = layout
    py.iplot ( fig )
```

## STATISTICAL CHARTS

### Histograms

```
trace = go.Histogram(  
    x = [ 1, 2, 3, 3, 3, 4, 5 ])  
data = [ trace ]  
py.iplot ( data )
```

### Box Plots

```
trace = go.Box (  
    x = [ 1, 2, 3, 3, 3, 4, 5 ])  
data = [ trace ]  
py.iplot ( data )
```

### 2D Histogram

```
trace = go.Histogram2d (  
    x = [ 1, 2, 3, 3, 3, 4, 5 ],  
    y = [ 1, 2, 3, 3, 3, 4, 5 ])  
data = [ trace ]  
py.iplot ( data )
```

## MAPS

### Bubble Map

```
trace = dict (  
    type = 'scattergeo',  
    lon = [ 100, 400 ], lat = [ 0, 0 ],  
    marker = dict (  
        marker = [ 'red', 'blue' ]  
        size = [ 30, 50 ] ),  
    mode = 'markers' )  
py.iplot ([ trace ])
```

### Choropleth Map

```
trc = dict (  
    type = 'choropleth',  
    locations = [ 'AZ', 'CA', 'VT' ],  
    locationmode = 'USA-states',  
    colorscale = [ 'Viridis' ],  
    z = [ 10, 20, 40 ] )  
lyt = dict ( geo = dict ( scope = 'usa' ) )  
map = go.Figure ( data = [ trc ],  
    layout = lyt )  
py.iplot ( map )
```

### Scatter Map

```
trace = dict (  
    type = 'scattergeo',  
    lon = [ 42, 39 ], lat = [ 12, 22 ],  
    marker = [ 'Rome', 'Greece' ],  
    mode = 'markers' )  
py.iplot ([ trace ])
```

## 3D CHARTS

### 3D Surface Plots

```
trace = go.Surface (  
    colorscale = 'Viridis',  
    z = [ [ 3, 5, 8, 13 ],  
          [ 21, 13, 8, 5 ] ] )  
data = [ trace ]  
py.iplot ( data )
```

### 3D Line Plots

```
trace = go.Scatter3D (  
    x = [ 9, 8, 5, 1 ], y = [ 1, 2, 4, 8 ],  
    z = [ 11, 8, 15, 3 ],  
    mode = 'lines' )  
data = [ trace ]  
py.iplot ( data )
```

### 3D Scatter Plots

```
trace = go.Scatter3D (  
    x = [ 9, 8, 5, 1 ], y = [ 1, 2, 4, 8 ],  
    z = [ 11, 8, 15, 3 ],  
    mode = 'markers' )  
data = [ trace ]  
py.iplot ( data )
```

## FIGURE HIERARCHY

### Figure {}

DATA []  
TRACE {}  
x, y, z []  
color, text, size []  
colorscale ABC or []  
MARKER {}  
color ABC  
symbol ABC  
LINE {}  
color ABC  
width 123

LAYOUT {}  
title ABC  
XAXIS, YAXIS {}  
SCENE {}  
XAXIS, YAXIS, ZAXIS {}  
GEO {}  
LEGEND {}  
ANNOTATIONS {}

{ } = dictionary  
[ ] = list  
ABC = string  
123 = number

# Python For Data Science Cheat Sheet

## Bokeh

Learn Bokeh [Interactively](#) at [www.DataCamp.com](http://www.DataCamp.com), taught by Bryan Van de Ven, core contributor

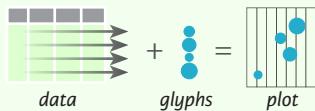


### Plotting With Bokeh

The Python interactive visualization library **Bokeh** enables high-performance visual presentation of large datasets in modern web browsers.



Bokeh's mid-level general purpose `bokeh.plotting` interface is centered around two main components: data and glyphs.



The basic steps to creating plots with the `bokeh.plotting` interface are:

1. Prepare some data:  
Python lists, NumPy arrays, Pandas DataFrames and other sequences of values
2. Create a new plot
3. Add renderers for your data, with visual customizations
4. Specify where to generate the output
5. Show or save the results

```
>>> from bokeh.plotting import figure
>>> from bokeh.io import output_file, show
>>> x = [1, 2, 3, 4, 5]           Step 1
>>> y = [6, 7, 2, 4, 5]
>>> p = figure(title="simple line example",      Step 2
              x_axis_label='x',
              y_axis_label='y')
>>> p.line(x, y, legend="Temp.", line_width=2)    Step 3
>>> output_file("lines.html")                   Step 4
>>> show(p)                                     Step 5
```

## 1 Data

### Also see Lists, NumPy & Pandas

Under the hood, your data is converted to Column Data Sources. You can also do this manually:

```
>>> import numpy as np
>>> import pandas as pd
>>> df = pd.DataFrame(np.array([[33.9, 4, 65, 'US'],
                               [32.4, 4, 66, 'Asia'],
                               [21.4, 4, 109, 'Europe']]),
                     columns=['mpg', 'cyl', 'hp', 'origin'],
                     index=['Toyota', 'Fiat', 'Volvo'])

>>> from bokeh.models import ColumnDataSource
>>> cds_df = ColumnDataSource(df)
```

## 2 Plotting

```
>>> from bokeh.plotting import figure
>>> p1 = figure(plot_width=300, tools='pan,box_zoom')
>>> p2 = figure(plot_width=300, plot_height=300,
               x_range=(0, 8), y_range=(0, 8))
>>> p3 = figure()
```

## 3 Renderers & Visual Customizations

### Glyphs

#### Scatter Markers

```
>>> p1.circle(np.array([1,2,3]), np.array([3,2,1]),
             fill_color='white')
>>> p2.square(np.array([1.5,3.5,5.5]), [1,4,3],
             color='blue', size=1)
```

#### Line Glyphs

```
>>> p1.line([1,2,3,4], [3,4,5,6], line_width=2)
>>> p2.multi_line(pd.DataFrame([[1,2,3],[5,6,7]]),
                  pd.DataFrame([[3,4,5],[3,2,1]]),
                  color="blue")
```

### Rows & Columns Layout

#### Rows

```
>>> from bokeh.layouts import row
```

```
>>> layout = row(p1,p2,p3)
```

#### Nesting Rows & Columns

```
>>> layout = row(column(p1,p2), p3)
```

#### Columns

```
>>> from bokeh.layouts import column
```

```
>>> layout = column(p1,p2,p3)
```

### Grid Layout

```
>>> from bokeh.layouts import gridplot
>>> row1 = [p1,p2]
>>> row2 = [p3]
>>> layout = gridplot([[p1,p2], [p3]])
```

### Tabbed Layout

```
>>> from bokeh.models.widgets import Panel, Tabs
>>> tab1 = Panel(child=p1, title="tab1")
>>> tab2 = Panel(child=p2, title="tab2")
>>> layout = Tabs(tabs=[tab1, tab2])
```

### Legends

#### Legend Location

##### Inside Plot Area

```
>>> p.legend.location = 'bottom_left'
```

##### Outside Plot Area

```
>>> r1 = p2.asterisk(np.array([1,2,3]), np.array([3,2,1]))
>>> r2 = p2.line([1,2,3,4], [3,4,5,6])
>>> legend = Legend(items=[("One", [p1, r1]), ("Two", [r2])], location=(0, -30))
>>> p.add_layout(legend, 'right')
```

### Customized Glyphs

#### Selection and Non-Selection Glyphs

```
>>> p = figure(tools='box_select')
>>> p.circle('mpg', 'cyl', source=cds_df,
             selection_color='red',
             nonselection_alpha=0.1)
```

#### Hover Glyphs

```
>>> hover = HoverTool(tooltips=None, mode='vline')
>>> p3.add_tools(hover)
```

#### Colormapping

```
>>> color_mapper = CategoricalColorMapper(
            factors=['US', 'Asia', 'Europe'],
            palette=['blue', 'red', 'green'])
>>> p3.circle('mpg', 'cyl', source=cds_df,
             color=dict(field='origin',
                        transform=color_mapper),
             legend='Origin'))
```

### Also see Data

## 4 Output

### Output to HTML File

```
>>> from bokeh.io import output_file, show
>>> output_file('my_bar_chart.html', mode='cdn')
```

### Notebook Output

```
>>> from bokeh.io import output_notebook, show
>>> output_notebook()
```

### Embedding

#### Standalone HTML

```
>>> from bokeh.embed import file_html
>>> html = file_html(p, CDN, "my_plot")
```

#### Components

```
>>> from bokeh.embed import components
>>> script, div = components(p)
```

## 5 Show or Save Your Plots

```
>>> show(p1)
>>> show(layout)
```

```
>>> save(p1)
>>> save(layout)
```

### Also see Data

## Statistical Charts With Bokeh

### Also see Data

Bokeh's high-level `bokeh.charts` interface is ideal for quickly creating statistical charts

### Bar Chart

```
>>> from bokeh.charts import Bar
>>> p = Bar(df, stacked=True, palette=['red','blue'])
```

### Box Plot

```
>>> from bokeh.charts import BoxPlot
>>> p = BoxPlot(df, values='vals', label='cyl',
                 legend='bottom_right')
```

### Histogram

```
>>> from bokeh.charts import Histogram
>>> p = Histogram(df, title='Histogram')
```

### Scatter Plot

```
>>> from bokeh.charts import Scatter
>>> p = Scatter(df, x='mpg', y ='hp', marker='square',
                 xlabel='Miles Per Gallon',
                 ylabel='Horsepower')
```

