# Beginner's Python Cheat Sheet

## Variables and Strings

*Variables are used to store values. A string is a series of characters, surrounded by single or double quotes.*

### Hello world

```python
print("Hello world!")
```

### Hello world with a variable

```python
msg = "Hello world!"
print(msg)
```

### Concatenation (combining strings)

```python
first_name = 'albert'
last_name = 'einstein'
full_name = first_name + ' ' + last_name
print(full_name)
```

## Lists

*A list stores a series of items in a particular order. You access items using an index, or within a loop.*

### Make a list

```python
bikes = ['trek', 'redline', 'giant']
```

### Get the first item in a list

```python
first_bike = bikes[0]
```

### Get the last item in a list

```python
last_bike = bikes[-1]
```

### Looping through a list

```python
for bike in bikes:
    print(bike)
```

### Adding items to a list

```python
bikes = []
bikes.append('trek')
bikes.append('redline')
bikes.append('giant')
```

### Making numerical lists

```python
squares = []
for x in range(1, 11):
    squares.append(x**2)
```

## Lists (cont.)

### List comprehensions

```python
squares = [x**2 for x in range(1, 11)]
```

### Slicing a list

```python
finishers = ['sam', 'bob', 'ada', 'bea']
first_two = finishers[:2]
```

### Copying a list

```python
copy_of_bikes = bikes[:]
```

## Tuples

*Tuples are similar to lists, but the items in a tuple can't be modified.*

### Making a tuple

```python
dimensions = (1920, 1080)
```

## If statements

*If statements are used to test for particular conditions and respond appropriately.*

### Conditional tests

```
equals              x == 42
not equal           x != 42
greater than        x > 42
  or equal to       x >= 42
less than           x < 42
  or equal to       x <= 42
```

### Conditional test with lists

```python
'trek' in bikes
'surly' not in bikes
```

### Assigning boolean values

```python
game_active = True
can_edit = False
```

### A simple if test

```python
if age >= 18:
    print("You can vote!")
```

### If-elif-else statements

```python
if age < 4:
    ticket_price = 0
elif age < 18:
    ticket_price = 10
else:
    ticket_price = 15
```

## Dictionaries

*Dictionaries store connections between pieces of information. Each item in a dictionary is a key-value pair.*

### A simple dictionary

```python
alien = {'color': 'green', 'points': 5}
```

### Accessing a value

```python
print("The alien's color is " + alien['color'])
```

### Adding a new key-value pair

```python
alien['x_position'] = 0
```

### Looping through all key-value pairs

```python
fav_numbers = {'eric': 17, 'ever': 4}
for name, number in fav_numbers.items():
    print(name + ' loves ' + str(number))
```

### Looping through all keys

```python
fav_numbers = {'eric': 17, 'ever': 4}
for name in fav_numbers.keys():
    print(name + ' loves a number')
```

### Looping through all the values

```python
fav_numbers = {'eric': 17, 'ever': 4}
for number in fav_numbers.values():
    print(str(number) + ' is a favorite')
```

## User input

*Your programs can prompt the user for input. All input is stored as a string.*

### Prompting for a value

```python
name = input("What's your name? ")
print("Hello, " + name + "!")
```

### Prompting for numerical input

```python
age = input("How old are you? ")
age = int(age)

pi = input("What's the value of pi? ")
pi = float(pi)
```

## While loops

*A while loop repeats a block of code as long as a certain condition is true.*

### A simple while loop

```python
current_value = 1
while current_value <= 5:
    print(current_value)
    current_value += 1
```

### Letting the user choose when to quit

```python
msg = ''
while msg != 'quit':
    msg = input("What's your message? ")
    print(msg)
```

## Functions

*Functions are named blocks of code, designed to do one specific job. Information passed to a function is called an argument, and information received by a function is called a parameter.*

### A simple function

```python
def greet_user():
    """Display a simple greeting."""
    print("Hello!")

greet_user()
```

### Passing an argument

```python
def greet_user(username):
    """Display a personalized greeting."""
    print("Hello, " + username + "!")

greet_user('jesse')
```

### Default values for parameters

```python
def make_pizza(topping='bacon'):
    """Make a single-topping pizza."""
    print("Have a " + topping + " pizza!")

make_pizza()
make_pizza('pepperoni')
```

### Returning a value

```python
def add_numbers(x, y):
    """Add two numbers and return the sum."""
    return x + y

sum = add_numbers(3, 5)
print(sum)
```

## Classes

*A class defines the behavior of an object and the kind of information an object can store. The information in a class is stored in attributes, and functions that belong to a class are called methods. A child class inherits the attributes and methods from its parent class.*

### Creating a dog class

```python
class Dog():
    """Represent a dog."""

    def __init__(self, name):
        """Initialize dog object."""
        self.name = name

    def sit(self):
        """Simulate sitting."""
        print(self.name + " is sitting.")

my_dog = Dog('Peso')

print(my_dog.name + " is a great dog!")
my_dog.sit()
```

### Inheritance

```python
class SARDog(Dog):
    """Represent a search dog."""

    def __init__(self, name):
        """Initialize the sardog."""
        super().__init__(name)

    def search(self):
        """Simulate searching."""
        print(self.name + " is searching.")

my_dog = SARDog('Willie')

print(my_dog.name + " is a search dog.")
my_dog.sit()
my_dog.search()
```

## Working with files

*Your programs can read from files and write to files. Files are opened in read mode ('r') by default, but can also be opened in write mode ('w') and append mode ('a').*

### Reading a file and storing its lines

```python
filename = 'siddhartha.txt'
with open(filename) as file_object:
    lines = file_object.readlines()

for line in lines:
    print(line)
```

### Writing to a file

```python
filename = 'journal.txt'
with open(filename, 'w') as file_object:
    file_object.write("I love programming.")
```

### Appending to a file

```python
filename = 'journal.txt'
with open(filename, 'a') as file_object:
    file_object.write("\nI love making games.")
```

## Exceptions

*Exceptions help you respond appropriately to errors that are likely to occur. You place code that might cause an error in the try block. Code that should run in response to an error goes in the except block. Code that should run only if the try block was successful goes in the else block.*

### Catching an exception

```python
prompt = "How many tickets do you need? "
num_tickets = input(prompt)

try:
    num_tickets = int(num_tickets)
except ValueError:
    print("Please try again.")
else:
    print("Your tickets are printing.")
```

# Python For Data Science *Cheat Sheet*
## Python Basics

## Variables and Data Types

### Variable Assignment
```
>>> x=5
>>> x
 5
```

### Calculations With Variables

| | |
|---|---|
| `>>> x+2`<br>` 7` | Sum of two variables |
| `>>> x-2`<br>` 3` | Subtraction of two variables |
| `>>> x*2`<br>` 10` | Multiplication of two variables |
| `>>> x**2`<br>` 25` | Exponentiation of a variable |
| `>>> x%2`<br>` 1` | Remainder of a variable |
| `>>> x/float(2)`<br>` 2.5` | Division of a variable |

### Types and Type Conversion

| | | |
|---|---|---|
| `str()` | `'5'`, `'3.45'`, `'True'` | Variables to strings |
| `int()` | `5`, `3`, `1` | Variables to integers |
| `float()` | `5.0`, `1.0` | Variables to floats |
| `bool()` | `True`, `True`, `True` | Variables to booleans |

## Asking For Help
```
>>> help(str)
```

## Strings
```
>>> my_string = 'thisStringIsAwesome'
>>> my_string
'thisStringIsAwesome'
```

### String Operations
```
>>> my_string * 2
 'thisStringIsAwesomethisStringIsAwesome'
>>> my_string + 'Innit'
 'thisStringIsAwesomeInnit'
>>> 'm' in my_string
 True
```

## Lists

**Also see NumPy Arrays**

```
>>> a = 'is'
>>> b = 'nice'
>>> my_list = ['my', 'list', a, b]
>>> my_list2 = [[4,5,6,7], [3,4,5,6]]
```

### Selecting List Elements

**Index starts at 0**

#### Subset
```
>>> my_list[1]          Select item at index 1
>>> my_list[-3]         Select 3rd last item
```
#### Slice
```
>>> my_list[1:3]        Select items at index 1 and 2
>>> my_list[1:]         Select items after index 0
>>> my_list[:3]         Select items before index 3
>>> my_list[:]          Copy my_list
```
#### Subset Lists of Lists
```
>>> my_list2[1][0]      my_list[list][itemOfList]
>>> my_list2[1][:2]
```

### List Operations
```
>>> my_list + my_list
['my', 'list', 'is', 'nice', 'my', 'list', 'is', 'nice']
>>> my_list * 2
['my', 'list', 'is', 'nice', 'my', 'list', 'is', 'nice']
>>> my_list2 > 4
True
```

### List Methods

| | |
|---|---|
| `>>> my_list.index(a)` | Get the index of an item |
| `>>> my_list.count(a)` | Count an item |
| `>>> my_list.append('!')` | Append an item at a time |
| `>>> my_list.remove('!')` | Remove an item |
| `>>> del(my_list[0:1])` | Remove an item |
| `>>> my_list.reverse()` | Reverse the list |
| `>>> my_list.extend('!')` | Append an item |
| `>>> my_list.pop(-1)` | Remove an item |
| `>>> my_list.insert(0,'!')` | Insert an item |
| `>>> my_list.sort()` | Sort the list |

### String Operations

**Index starts at 0**

```
>>> my_string[3]
>>> my_string[4:9]
```

### String Methods

| | |
|---|---|
| `>>> my_string.upper()` | String to uppercase |
| `>>> my_string.lower()` | String to lowercase |
| `>>> my_string.count('w')` | Count String elements |
| `>>> my_string.replace('e', 'i')` | Replace String elements |
| `>>> my_string.strip()` | Strip whitespaces |

## Libraries

### Import libraries
```
>>> import numpy
>>> import numpy as np
```
*pandas* — Data analysis
*learn* — Machine learning

### Selective import
```
>>> from math import pi
```
*NumPy* — Scientific computing
*matplotlib* — 2D plotting

## Install Python

**ANACONDA**
Leading open data science platform powered by Python

**spyder**
Free IDE that is included with Anaconda

**jupyter**
Create and share documents with live code, visualizations, text, ...

## Numpy Arrays

**Also see Lists**

```
>>> my_list = [1, 2, 3, 4]
>>> my_array = np.array(my_list)
>>> my_2darray = np.array([[1,2,3],[4,5,6]])
```

### Selecting Numpy Array Elements

**Index starts at 0**

#### Subset
```
>>> my_array[1]         Select item at index 1
 2
```
#### Slice
```
>>> my_array[0:2]       Select items at index 0 and 1
 array([1, 2])
```
#### Subset 2D Numpy arrays
```
>>> my_2darray[:,0]     my_2darray[rows, columns]
 array([1, 4])
```

### Numpy Array Operations
```
>>> my_array > 3
 array([False, False, False,  True], dtype=bool)
>>> my_array * 2
 array([2, 4, 6, 8])
>>> my_array + np.array([5, 6, 7, 8])
 array([6, 8, 10, 12])
```

### Numpy Array Functions

| | |
|---|---|
| `>>> my_array.shape` | Get the dimensions of the array |
| `>>> np.append(other_array)` | Append items to an array |
| `>>> np.insert(my_array, 1, 5)` | Insert items in an array |
| `>>> np.delete(my_array, [1])` | Delete items in an array |
| `>>> np.mean(my_array)` | Mean of the array |
| `>>> np.median(my_array)` | Median of the array |
| `>>> my_array.corrcoef()` | Correlation coefficient |
| `>>> np.std(my_array)` | Standard deviation |

# Python For Data Science *Cheat Sheet*
## Importing Data

## 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')    Open the file for reading
>>> text = file.read()                 Read a file's contents
>>> print(file.closed)                 Check whether file is closed
>>> file.close()                       Close file
>>> print(text)
```

#### Using the context manager `with`

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

### Table Data: Flat Files

#### Importing Flat Files with numpy

**Files with one data type**

```
>>> filename = 'mnist.txt'
>>> data = np.loadtxt(filename,
                      delimiter=',',    String used to separate values
                      skiprows=2,       Skip the first 2 lines
                      usecols=[0,2],    Read the 1st and 3rd column
                      dtype=str)        The type of the resulting array
```

**Files with mixed data types**

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

```
>>> 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,          Number of rows of file to read
                       header=None,      Row number to use as col names
                       sep='\t',         Delimiter to use
                       comment='#',      Character to split comments
                       na_values=[""])   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.sas7bdat') 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 type of array elements
>>> data_array.shape     Array dimensions
>>> len(data_array)      Length of array
```

### pandas DataFrames

```
>>> df.head()                        Return first DataFrame rows
>>> df.tail()                        Return last DataFrame rows
>>> df.index                         Describe index
>>> df.columns                       Describe DataFrame columns
>>> df.info()                        Info on DataFrame
>>> data_array = data.values         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

```
>>> print(mat.keys())          Print dictionary keys
>>> for key in data.keys():    Print dictionary keys
        print(key)
meta
quality
strain
>>> pickled_data.values()      Return dictionary values
>>> print(mat.items())         Returns items in list format of (key, value)
                               tuple pairs
```

### Accessing Data Items with Keys

```
>>> for key in data ['meta'].keys()    Explore the HDF5 structure
        print(key)
Description
DescriptionURL
Detector
Duration
GPSstart
Observatory
Type
UTCstart
>>> print(data['meta']['Description'].value)    Retrieve the value for a key
```

## Navigating Your FileSystem

### Magic Commands

```
!ls      List directory contents of files and directories
%cd ..   Change current working directory
%pwd     Return the current working directory path
```

### `os` Library

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

# Python For Data Science *Cheat Sheet*
## Pandas Basics

## Pandas

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

$$pandas$$
$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$

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

Columns →

| | Country | Capital | Population |
|---|---|---|---|
| 0 | Belgium | Brussels | 11190846 |
| 1 | India | New Delhi | 1303171035 |
| 2 | Brazil | Brasília | 207847528 |

Index →

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

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],
            'Capital': ['Brussels', 'New Delhi', 'Brasília'],
            '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                                   **Also see NumPy Arrays**

### Getting

| | |
|---|---|
| ```>>> s['b']```<br>`  -5` | Get one element |
| ```>>> df[1:]```<br>`   Country    Capital   Population`<br>`1   India   New Delhi   1303171035`<br>`2   Brazil   Brasília   207847528` | Get subset of a DataFrame |

### Selecting, Boolean Indexing & Setting

#### By Position

| | |
|---|---|
| ```>>> df.iloc[[0],[0]]```<br>`'Belgium'`<br>```>>> df.iat([0],[0])```<br>`'Belgium'` | Select single value by row & column |

#### By Label

| | |
|---|---|
| ```>>> df.loc[[0], ['Country']]```<br>`'Belgium'`<br>```>>> df.at([0], ['Country'])```<br>`'Belgium'` | Select single value by row & column labels |

#### By Label/Position

| | |
|---|---|
| ```>>> df.ix[2]```<br>`Country    Brazil`<br>`Capital    Brasília`<br>`Population    207847528` | Select single row of subset of rows |
| ```>>> df.ix[:,'Capital']```<br>`0    Brussels`<br>`1    New Delhi`<br>`2    Brasília` | Select a single column of subset of columns |
| ```>>> df.ix[1,'Capital']```<br>`'New Delhi'` | Select rows and columns |

#### Boolean Indexing

| | |
|---|---|
| ```>>> s[~(s > 1)]```<br>```>>> s[(s < -1) | (s > 2)]```<br>```>>> df[df['Population']>1200000000]``` | Series `s` where value is not >1<br>`s`  where value is <-1 or >2<br>Use filter to adjust DataFrame |

#### Setting

| | |
|---|---|
| ```>>> s['a'] = 6``` | Set index `a` of Series `s` to 6 |

### Read and Write to SQL Query or Database Table

```
>>> from sqlalchemy import create_engine
>>> engine = create_engine('sqlite:///:memory:')
>>> pd.read_sql("SELECT * FROM my_table;", engine)
>>> pd.read_sql_table('my_table', engine)
>>> pd.read_sql_query("SELECT * FROM my_table;", engine)
```
`read_sql()` is a convenience wrapper around `read_sql_table()` and
`read_sql_query()`

```
>>> pd.to_sql('myDf', engine)
```

## Dropping

| | |
|---|---|
| ```>>> s.drop(['a', 'c'])``` | Drop values from rows (axis=0) |
| ```>>> df.drop('Country', axis=1)``` | Drop values from columns(axis=1) |

## Sort & Rank

| | |
|---|---|
| ```>>> df.sort_index()``` | Sort by labels along an axis |
| ```>>> df.sort_values(by='Country')``` | Sort by the values along an axis |
| ```>>> df.rank()``` | Assign ranks to entries |

## Retrieving Series/DataFrame Information

### Basic Information

| | |
|---|---|
| ```>>> df.shape``` | (rows,columns) |
| ```>>> df.index``` | Describe index |
| ```>>> df.columns``` | Describe DataFrame columns |
| ```>>> df.info()``` | Info on DataFrame |
| ```>>> df.count()``` | Number of non-NA values |

### Summary

| | |
|---|---|
| ```>>> df.sum()``` | Sum of values |
| ```>>> df.cumsum()``` | Cummulative sum of values |
| ```>>> df.min()/df.max()``` | Minimum/maximum values |
| ```>>> df.idxmin()/df.idxmax()``` | Minimum/Maximum index value |
| ```>>> df.describe()``` | Summary statistics |
| ```>>> df.mean()``` | Mean of values |
| ```>>> df.median()``` | Median of values |

## Applying Functions

| | |
|---|---|
| ```>>> f = lambda x: x*2```<br>```>>> df.apply(f)``` | Apply function |
| ```>>> df.applymap(f)``` | 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)
```

# Python For Data Science *Cheat Sheet*
## Pandas

## Reshaping Data

### Pivot

```
>>> df3= df2.pivot(index='Date',
                   columns='Type',
                   values='Value')
```
Spread rows into columns

| | Date | Type | Value |
|---|---|---|---|
| 0 | 2016-03-01 | a | 11.432 |
| 1 | 2016-03-02 | b | 13.031 |
| 2 | 2016-03-01 | c | 20.784 |
| 3 | 2016-03-03 | a | 99.906 |
| 4 | 2016-03-02 | a | 1.303 |
| 5 | 2016-03-03 | c | 20.784 |

| Type | a | b | c |
|---|---|---|---|
| Date | | | |
| 2016-03-01 | 11.432 | NaN | 20.784 |
| 2016-03-02 | 1.303 | 13.031 | NaN |
| 2016-03-03 | 99.906 | NaN | 20.784 |

### Pivot Table

```
>>> df4 = pd.pivot_table(df2,
                   values='Value',
                   index='Date',
                   columns='Type'])
```
Spread rows into columns

### Stack / Unstack

```
>>> stacked = df5.stack()
>>> stacked.unstack()
```
Pivot a level of column labels
Pivot a level of index labels

| | | 0 | 1 |
|---|---|---|---|
| 1 | 5 | 0.233482 | 0.390959 |
| 2 | 4 | 0.184713 | 0.237102 |
| 3 | 3 | 0.433522 | 0.429401 |

*Unstacked*

| | | | |
|---|---|---|---|
| 1 | 5 | 0 | 0.233482 |
| | | 1 | 0.390959 |
| 2 | 4 | 0 | 0.184713 |
| | | 1 | 0.237102 |
| 3 | 3 | 0 | 0.433522 |
| | | 1 | 0.429401 |

*Stacked*

### Melt

```
>>> pd.melt(df2,
           id_vars=["Date"],
           value_vars=["Type", "Value"],
           value_name="Observations")
```
Gather columns into rows

| | Date | Type | Value |
|---|---|---|---|
| 0 | 2016-03-01 | a | 11.432 |
| 1 | 2016-03-02 | b | 13.031 |
| 2 | 2016-03-01 | c | 20.784 |
| 3 | 2016-03-03 | a | 99.906 |
| 4 | 2016-03-02 | a | 1.303 |
| 5 | 2016-03-03 | c | 20.784 |

| | Date | Variable | Observations |
|---|---|---|---|
| 0 | 2016-03-01 | Type | a |
| 1 | 2016-03-02 | Type | b |
| 2 | 2016-03-01 | Type | c |
| 3 | 2016-03-03 | Type | a |
| 4 | 2016-03-02 | Type | a |
| 5 | 2016-03-03 | Type | c |
| 6 | 2016-03-01 | Value | 11.432 |
| 7 | 2016-03-02 | Value | 13.031 |
| 8 | 2016-03-01 | Value | 20.784 |
| 9 | 2016-03-03 | Value | 99.906 |
| 10 | 2016-03-02 | Value | 1.303 |
| 11 | 2016-03-03 | Value | 20.784 |

## Iteration

```
>>> df.iteritems()
>>> df.iterrows()
```
(Column-index, Series) pairs
(Row-index, Series) pairs

## Advanced Indexing
**Also see NumPy Arrays**

### Selecting
```
>>> df3.loc[:,(df3>1).any()]
>>> df3.loc[:,(df3>1).all()]
>>> df3.loc[:,df3.isnull().any()]
>>> df3.loc[:,df3.notnull().all()]
```
Select cols with any vals >1
Select cols with vals > 1
Select cols with NaN
Select cols without NaN

### Indexing With isin
```
>>> df[(df.Country.isin(df2.Type))]
>>> df3.filter(items="a","b")]
>>> df.select(lambda x: not x%5)
```
Find same elements
Filter on values
Select specific elements

### Where
```
>>> s.where(s > 0)
```
Subset the data

### Query
```
>>> df6.query('second > first')
```
Query DataFrame

### Setting/Resetting Index
```
>>> df.set_index('Country')
>>> df4 = df.reset_index()
>>> df = df.rename(index=str,
                columns={"Country":"cntry",
                         "Capital":"cptl",
                         "Population":"ppltn"})
```
Set the index
Reset the index
Rename DataFrame

### Reindexing
```
>>> s2 = s.reindex(['a','c','d','e','b'])
```

#### Forward Filling
```
>>> df.reindex(range(4),
            method='ffill')
   Country     Capital   Population
0  Belgium    Brussels   11190846
1  India      New Delhi  1303171035
2  Brazil     Brasília   207847528
3  Brazil     Brasília   207847528
```

#### Backward Filling
```
>>> s3 = s.reindex(range(5),
            method='bfill')
0    3
1    3
2    3
3    3
4    3
```

### MultiIndexing
```
>>> arrays = [np.array([1,2,3]),
              np.array([5,4,3])]
>>> df5 = pd.DataFrame(np.random.rand(3, 2), index=arrays)
>>> tuples = list(zip(*arrays))
>>> index = pd.MultiIndex.from_tuples(tuples,
                          names=['first', 'second'])
>>> df6 = pd.DataFrame(np.random.rand(3, 2), index=index)
>>> df2.set_index(["Date", "Type"])
```

### Duplicate Data
```
>>> s3.unique()
>>> df2.duplicated('Type')
>>> df2.drop_duplicates('Type', keep='last')
>>> df.index.duplicated()
```
Return unique values
Check duplicates
Drop duplicates
Check index duplicates

### Grouping Data

#### Aggregation
```
>>> df2.groupby(by=['Date','Type']).mean()
>>> df4.groupby(level=0).sum()
>>> df4.groupby(level=0).agg({'a':lambda x:sum(x)/len(x),
                              'b': np.sum})
```
#### Transformation
```
>>> customSum = lambda x: (x+x%2)
>>> df4.groupby(level=0).transform(customSum)
```

### Missing Data
```
>>> df.dropna()
>>> df3.fillna(df3.mean())
>>> df2.replace("a", "f")
```
Drop NaN values
Fill NaN values with a predetermined value
Replace values with others

## Combining Data

*data1*

| X1 | X2 |
|---|---|
| a | 11.432 |
| b | 1.303 |
| c | 99.906 |

*data2*

| X1 | X3 |
|---|---|
| a | 20.784 |
| b | NaN |
| d | 20.784 |

### Merge
```
>>> pd.merge(data1,
             data2,
             how='left',
             on='X1')
```
| X1 | X2 | X3 |
|---|---|---|
| a | 11.432 | 20.784 |
| b | 1.303 | NaN |
| c | 99.906 | NaN |

```
>>> pd.merge(data1,
             data2,
             how='right',
             on='X1')
```
| X1 | X2 | X3 |
|---|---|---|
| a | 11.432 | 20.784 |
| b | 1.303 | NaN |
| d | NaN | 20.784 |

```
>>> pd.merge(data1,
             data2,
             how='inner',
             on='X1')
```
| X1 | X2 | X3 |
|---|---|---|
| a | 11.432 | 20.784 |
| b | 1.303 | NaN |

```
>>> pd.merge(data1,
             data2,
             how='outer',
             on='X1')
```
| X1 | X2 | X3 |
|---|---|---|
| a | 11.432 | 20.784 |
| b | 1.303 | NaN |
| c | 99.906 | NaN |
| d | NaN | 20.784 |

### Join
```
>>> data1.join(data2, how='right')
```

### Concatenate
#### Vertical
```
>>> s.append(s2)
```
#### Horizontal/Vertical
```
>>> pd.concat([s,s2],axis=1, keys=['One','Two'])
>>> pd.concat([data1, data2], axis=1, join='inner')
```

### Dates
```
>>> df2['Date']= pd.to_datetime(df2['Date'])
>>> df2['Date']= pd.date_range('2000-1-1',
                    periods=6,
                    freq='M')
>>> dates = [datetime(2012,5,1), datetime(2012,5,2)]
>>> index = pd.DatetimeIndex(dates)
>>> index = pd.date_range(datetime(2012,2,1), end, freq='BM')
```

### Visualization
**Also see Matplotlib**
```
>>> import matplotlib.pyplot as plt
```
```
>>> s.plot()
>>> plt.show()
```
```
>>> df2.plot()
>>> plt.show()
```

# Python For Data Science *Cheat Sheet*
## NumPy Basics

## 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.
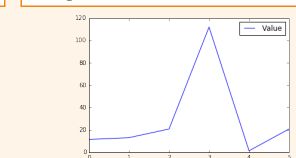
Use the following import convention:

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

### NumPy Arrays

**1D array**

| 1 | 2 | 3 |
|---|---|---|

**2D array**

axis 1
axis 0

| 1.5 | 2 | 3 |
|-----|---|---|
| 4 | 5 | 6 |

**3D array**

axis 2
axis 1
axis 0

## 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))` | Create an array of zeros |
| `>>> np.ones((2,3,4),dtype=np.int16)` | Create an array of ones |
| `>>> d = np.arange(10,25,5)` | Create an array of evenly spaced values (step value) |
| `>>> np.linspace(0,2,9)` | Create an array of evenly spaced values (number of samples) |
| `>>> e = np.full((2,2),7)` | Create a constant array |
| `>>> f = np.eye(2)` | Create a 2X2 identity matrix |
| `>>> np.random.random((2,2))` | Create an array with random values |
| `>>> np.empty((3,2))` | Create an empty array |

## I/O

### Saving & Loading On Disk

```
>>> np.save('my_array', a)
>>> np.savez('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

| | |
|---|---|
| `>>> np.int64` | Signed 64-bit integer types |
| `>>> np.float32` | Standard double-precision floating point |
| `>>> np.complex` | Complex numbers represented by 128 floats |
| `>>> np.bool` | Boolean type storing `TRUE` and `FALSE` values |
| `>>> np.object` | Python object type |
| `>>> np.string_` | Fixed-length string type |
| `>>> np.unicode_` | Fixed-length unicode type |

## Inspecting Your Array

| | |
|---|---|
| `>>> a.shape` | Array dimensions |
| `>>> len(a)` | Length of array |
| `>>> b.ndim` | Number of array dimensions |
| `>>> e.size` | Number of array elements |
| `>>> b.dtype` | Data type of array elements |
| `>>> b.dtype.name` | Name of data type |
| `>>> b.astype(int)` | Convert an array to a different type |

## Asking For Help

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

## Array Mathematics

### Arithmetic Operations

| | |
|---|---|
| `>>> g = a - b`<br>`  array([[-0.5, 0. , 0. ],`<br>`       [-3. , -3. , -3. ]])` | Subtraction |
| `>>> np.subtract(a,b)` | Subtraction |
| `>>> b + a`<br>`  array([[ 2.5, 4. , 6. ],`<br>`       [ 5. , 7. , 9. ]])` | Addition |
| `>>> np.add(b,a)` | Addition |
| `>>> a / b`<br>`  array([[ 0.66666667, 1.    , 1.    ],`<br>`       [ 0.25  , 0.4   , 0.5   ]])` | Division |
| `>>> np.divide(a,b)` | Division |
| `>>> a * b`<br>`  array([[ 1.5, 4. , 9. ],`<br>`       [ 4. , 10. , 18. ]])` | Multiplication |
| `>>> np.multiply(a,b)` | Multiplication |
| `>>> np.exp(b)` | Exponentiation |
| `>>> np.sqrt(b)` | Square root |
| `>>> np.sin(a)` | Print sines of an array |
| `>>> np.cos(b)` | Element-wise cosine |
| `>>> np.log(a)` | Element-wise natural logarithm |
| `>>> e.dot(f)`<br>`  array([[ 7., 7.],`<br>`       [ 7., 7.]])` | Dot product |

### Comparison

| | |
|---|---|
| `>>> a == b`<br>`  array([[False, True, True],`<br>`       [False, False, False]], dtype=bool)` | Element-wise comparison |
| `>>> a < 2`<br>`  array([True, False, False], dtype=bool)` | Element-wise comparison |
| `>>> np.array_equal(a, b)` | Array-wise comparison |

### Aggregate Functions

| | |
|---|---|
| `>>> a.sum()` | Array-wise sum |
| `>>> a.min()` | Array-wise minimum value |
| `>>> b.max(axis=0)` | Maximum value of an array row |
| `>>> b.cumsum(axis=1)` | Cumulative sum of the elements |
| `>>> a.mean()` | Mean |
| `>>> b.median()` | Median |
| `>>> a.corrcoef()` | Correlation coefficient |
| `>>> np.std(b)` | Standard deviation |

## Copying Arrays

| | |
|---|---|
| `>>> h = a.view()` | Create a view of the array with the same data |
| `>>> np.copy(a)` | Create a copy of the array |
| `>>> h = a.copy()` | Create a deep copy of the array |

## Sorting Arrays

| | |
|---|---|
| `>>> a.sort()` | Sort an array |
| `>>> c.sort(axis=0)` | Sort the elements of an array's axis |

## Subsetting, Slicing, Indexing

### Subsetting

| | |
|---|---|
| `>>> a[2]`<br>`  3` | Select the element at the 2nd index |
| `>>> b[1,2]`<br>`  6.0` | Select the element at row 1 column 2 (equivalent to `b[1][2]`) |

### Slicing

| | |
|---|---|
| `>>> a[0:2]`<br>`  array([1, 2])` | Select items at index 0 and 1 |
| `>>> b[0:2,1]`<br>`  array([ 2.,  5.])` | Select items at rows 0 and 1 in column 1 |
| `>>> b[:1]`<br>`  array([[1.5, 2., 3.]])` | Select all items at row 0 (equivalent to `b[0:1, :]`) |
| `>>> c[1,...]`<br>`  array([[[ 3., 2., 1.],`<br>`       [ 4., 5., 6.]]])` | Same as `[1,:,:]` |
| `>>> a[ : :-1]`<br>`  array([3, 2, 1])` | Reversed array `a` |

### Boolean Indexing

| | |
|---|---|
| `>>> a[a<2]`<br>`  array([1])` | Select elements from `a` less than 2 |

### Fancy Indexing

| | |
|---|---|
| `>>> b[[1, 0, 1, 0],[0, 1, 2, 0]]`<br>`  array([ 4., 2., 6., 1.5])` | Select elements (1,0),(0,1),(1,2) and (0,0) |
| `>>> b[[1, 0, 1, 0]][:,[0,1,2,0]]`<br>`  array([[ 4.,5., 6., 4.],`<br>`       [ 1.5,2., 3., 1.5],`<br>`       [ 4., 5., 6., 4.],`<br>`       [ 1.5,2., 3., 1.5]])` | Select a subset of the matrix's rows and columns |

## Array Manipulation

### Transposing Array

| | |
|---|---|
| `>>> i = np.transpose(b)` | Permute array dimensions |
| `>>> i.T` | Permute array dimensions |

### Changing Array Shape

| | |
|---|---|
| `>>> b.ravel()` | Flatten the array |
| `>>> g.reshape(3,-2)` | Reshape, but don't change data |

### Adding/Removing Elements

| | |
|---|---|
| `>>> h.resize((2,6))` | Return a new array with shape (2,6) |
| `>>> np.append(h,g)` | Append items to an array |
| `>>> np.insert(a, 1, 5)` | Insert items in an array |
| `>>> np.delete(a,[1])` | Delete items from an array |

### Combining Arrays

| | |
|---|---|
| `>>> np.concatenate((a,d),axis=0)`<br>`  array([ 1, 2, 3, 10, 15, 20])` | Concatenate arrays |
| `>>> np.vstack((a,b))`<br>`  array([[ 1. , 2. , 3. ],`<br>`       [ 1.5, 2. , 3. ],`<br>`       [ 4. , 5. , 6. ]])` | Stack arrays vertically (row-wise) |
| `>>> np.r_[e,f]` | Stack arrays vertically (row-wise) |
| `>>> np.hstack((e,f))`<br>`  array([[ 7., 7., 1., 0.],`<br>`       [ 7., 7., 0., 1.]])` | Stack arrays horizontally (column-wise) |
| `>>> np.column_stack((a,d))`<br>`  array([[ 1, 10],`<br>`       [ 2, 15],`<br>`       [ 3, 20]])` | Create stacked column-wise arrays |
| `>>> np.c_[a,d]` | Create stacked column-wise arrays |

### Splitting Arrays

| | |
|---|---|
| `>>> np.hsplit(a,3)`<br>`  [array([1]),array([2]),array([3])]` | Split the array horizontally at the 3rd index |
| `>>> np.vsplit(c,2)`<br>`  [array([[[ 1.5, 2. , 1. ],`<br>`       [ 2. , 15]]]),`<br>`  array([[[ 3., 2., 3.],`<br>`       [ 4., 5., 6.]]])]` | Split the array vertically at the 2nd index |

# Python For Data Science *Cheat Sheet*
## Matplotlib

## Matplotlib

**Matplotlib** is a Python 2D plotting library which produces publication-quality figures in a variety of hardcopy formats and interactive environments across platforms.

## Plot Anatomy & Workflow

### Plot Anatomy

### Workflow

The basic steps to creating plots with matplotlib are:

**1** Prepare data  **2** Create plot  **3** Plot  **4** Customize plot  **5** Save plot  **6** Show plot

```
>>> import matplotlib.pyplot as plt
>>> x = [1,2,3,4]          Step 1
>>> y = [10,20,25,30]
>>> fig = plt.figure()     Step 2
>>> ax = fig.add_subplot(111)  Step 3
>>> ax.plot(x, y, color='lightblue', linewidth=3)  Step 3, 4
>>> ax.scatter([2,4,6],
               [5,15,25],
               color='darkgreen',
               marker='^')
>>> ax.set_xlim(1, 6.5)
>>> plt.savefig('foo.png')
>>> plt.show()             Step 6
```

## 1 Prepare The Data

**Also see Lists & NumPy**

### 1D Data
```
>>> import numpy as np
>>> x = np.linspace(0, 10, 100)
>>> y = np.cos(x)
>>> z = np.sin(x)
```

### 2D Data or Images
```
>>> data = 2 * np.random.random((10, 10))
>>> data2 = 3 * np.random.random((10, 10))
>>> Y, X = np.mgrid[-3:3:100j, -3:3:100j]
>>> U = -1 - X**2 + Y
>>> V = 1 + X - Y**2
>>> from matplotlib.cbook import get_sample_data
>>> img = np.load(get_sample_data('axes_grid/bivariate_normal.npy'))
```

## 2 Create Plot
```
>>> import matplotlib.pyplot as plt
```

### Figure
```
>>> fig = plt.figure()
>>> fig2 = plt.figure(figsize=plt.figaspect(2.0))
```

### Axes

All plotting is done with respect to an `Axes`. In most cases, a subplot will fit your needs. A subplot is an axes on a grid system.
```
>>> fig.add_axes()
>>> ax1 = fig.add_subplot(221) # row-col-num
>>> ax3 = fig.add_subplot(212)
>>> fig3, axes = plt.subplots(nrows=2, ncols=2)
>>> fig4, axes2 = plt.subplots(ncols=3)
```

## 4 Customize Plot

### Colors, Color Bars & Color Maps
```
>>> plt.plot(x, x, x, x**2, x, x**3)
>>> ax.plot(x, y, alpha = 0.4)
>>> ax.plot(x, y, c='k')
>>> fig.colorbar(im, orientation='horizontal')
>>> im = ax.imshow(img,
                   cmap='seismic')
```

### Markers
```
>>> fig, ax = plt.subplots()
>>> ax.scatter(x,y,marker=".")
>>> ax.plot(x,y,marker="o")
```

### Linestyles
```
>>> plt.plot(x,y,linewidth=4.0)
>>> plt.plot(x,y,ls='solid')
>>> plt.plot(x,y,ls='--')
>>> plt.plot(x,y,'--',x**2,y**2,'-.')
>>> plt.setp(lines,color='r',linewidth=4.0)
```

### Text & Annotations
```
>>> ax.text(1,
            -2.1,
            'Example Graph',
            style='italic')
>>> ax.annotate("Sine",
                xy=(8, 0),
                xycoords='data',
                xytext=(10.5, 0),
                textcoords='data',
                arrowprops=dict(arrowstyle="->",
                                connectionstyle="arc3"),)
```

### Mathtext
```
>>> plt.title(r'$sigma_i=15$', fontsize=20)
```

### Limits, Legends & Layouts

**Limits & Autoscaling**
```
>>> ax.margins(x=0.0,y=0.1)         Add padding to a plot
>>> ax.axis('equal')               Set the aspect ratio of the plot to 1
>>> ax.set(xlim=[0,10.5],ylim=[-1.5,1.5])  Set limits for x-and y-axis
>>> ax.set_xlim(0,10.5)            Set limits for x-axis
```
**Legends**
```
>>> ax.set(title='An Example Axes',    Set a title and x-and y-axis labels
           ylabel='Y-Axis',
           xlabel='X-Axis')
>>> ax.legend(loc='best')          No overlapping plot elements
```
**Ticks**
```
>>> ax.xaxis.set(ticks=range(1,5),   Manually set x-ticks
                 ticklabels=[3,100,-12,"foo"])
>>> ax.tick_params(axis='y',         Make y-ticks longer and go in and out
                   direction='inout',
                   length=10)
```
**Subplot Spacing**
```
>>> fig3.subplots_adjust(wspace=0.5,   Adjust the spacing between subplots
                         hspace=0.3,
                         left=0.125,
                         right=0.9,
                         top=0.9,
                         bottom=0.1)
>>> fig.tight_layout()             Fit subplot(s) in to the figure area
```
**Axis Spines**
```
>>> ax1.spines['top'].set_visible(False)          Make the top axis line for a plot invisible
>>> ax1.spines['bottom'].set_position(('outward',10))  Move the bottom axis line outward
```

## 3 Plotting Routines

### 1D Data
```
>>> fig, ax = plt.subplots()
>>> lines = ax.plot(x,y)              Draw points with lines or markers connecting them
>>> ax.scatter(x,y)                   Draw unconnected points, scaled or colored
>>> axes[0,0].bar([1,2,3],[3,4,5])    Plot vertical rectangles (constant width)
>>> axes[1,0].barh([0.5,1,2.5],[0,1,2])  Plot horizontal rectangles (constant height)
>>> axes[1,1].axhline(0.45)           Draw a horizontal line across axes
>>> axes[0,1].axvline(0.65)           Draw a vertical line across axes
>>> ax.fill(x,y,color='blue')         Draw filled polygons
>>> ax.fill_between(x,y,color='yellow')  Fill between y-values and 0
```

### 2D Data or Images
```
>>> fig, ax = plt.subplots()
>>> im = ax.imshow(img,               Colormapped or RGB arrays
                   cmap='gist_earth',
                   interpolation='nearest',
                   vmin=-2,
                   vmax=2)
```

### Vector Fields
```
>>> axes[0,1].arrow(0,0,0.5,0.5)      Add an arrow to the axes
>>> axes[1,1].quiver(y,z)             Plot a 2D field of arrows
>>> axes[0,1].streamplot(X,Y,U,V)     Plot a 2D field of arrows
```

### Data Distributions
```
>>> ax1.hist(y)                       Plot a histogram
>>> ax3.boxplot(y)                    Make a box and whisker plot
>>> ax3.violinplot(z)                 Make a violin plot
```
```
>>> axes2[0].pcolor(data2)            Pseudocolor plot of 2D array
>>> axes2[0].pcolormesh(data)         Pseudocolor plot of 2D array
>>> CS = plt.contour(Y,X,U)           Plot contours
>>> axes2[2].contourf(data1)          Plot filled contours
>>> axes2[2]= ax.clabel(CS)           Label a contour plot
```

## 5 Save Plot

**Save figures**
```
>>> plt.savefig('foo.png')
```
**Save transparent figures**
```
>>> plt.savefig('foo.png', transparent=True)
```

## 6 Show Plot
```
>>> plt.show()
```

## Close & Clear
```
>>> plt.cla()        Clear an axis
>>> plt.clf()        Clear the entire figure
>>> plt.close()      Close a window
```

# Python For Data Science *Cheat Sheet*
## SciPy - Linear Algebra

## 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  <span>Also see NumPy</span>

```
>>> 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.vpslit(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(c)` | Return the real part of the array elements |
| `>>> np.imag(c)` | 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` | |
| `>>> np.unwrap(g)` | Unwrap |
| `>>> np.logspace(0,10,3)` | Create an array of evenly spaced values (log scale) |
| `>>> np.select([c<4],[c*2])` | Return values from a list of arrays depending on conditions |
| `>>> misc.factorial(a)` | Factorial |
| `>>> misc.comb(10,3,exact=True)` | Combine N things taken at k time |
| `>>> misc.central_diff_weights(3)` | Weights for Np-point central derivative |
| `>>> misc.derivative(myfunc,1.0)` | Find the n-th derivative of a function at a point |

## Linear Algebra  <span>Also see NumPy</span>

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` | Inverse |
| `>>> linalg.inv(A)` | Inverse |
| `>>> A.T` | Tranpose matrix |
| `>>> A.H` | Conjugate transposition |
| `>>> np.trace(A)` | Trace |

### Norm

| | |
|---|---|
| `>>> linalg.norm(A)` | Frobenius norm |
| `>>> linalg.norm(A,1)` | L1 norm (max column sum) |
| `>>> linalg.norm(A,np.inf)` | L inf norm (max row sum) |

### Rank

| | |
|---|---|
| `>>> np.linalg.matrix_rank(C)` | Matrix rank |

### Determinant

| | |
|---|---|
| `>>> linalg.det(A)` | Determinant |

### Solving linear problems

| | |
|---|---|
| `>>> linalg.solve(A,b)` | Solver for dense matrices |
| `>>> E = np.mat(a).T` | Solver for dense matrices |
| `>>> linalg.lstsq(D,E)` | Least-squares solution to linear matrix equation |

### Generalized inverse

| | |
|---|---|
| `>>> linalg.pinv(C)` | Compute the pseudo-inverse of a matrix (least-squares solver) |
| `>>> linalg.pinv2(C)` | Compute the pseudo-inverse of a matrix (SVD) |

## Creating Sparse Matrices

| | |
|---|---|
| `>>> F = np.eye(3, k=1)` | Create a 2X2 identity matrix |
| `>>> G = np.mat(np.identity(2))` | Create a 2x2 identity matrix |
| `>>> C[C > 0.5] = 0` | |
| `>>> H = sparse.csr_matrix(C)` | Compressed Sparse Row matrix |
| `>>> I = sparse.csc_matrix(D)` | Compressed Sparse Column matrix |
| `>>> J = sparse.dok_matrix(A)` | Dictionary Of Keys matrix |
| `>>> E.todense()` | Sparse matrix to full matrix |
| `>>> sparse.isspmatrix_csc(A)` | Identify sparse matrix |

## Sparse Matrix Routines

### Inverse

| | |
|---|---|
| `>>> sparse.linalg.inv(I)` | Inverse |

### Norm

| | |
|---|---|
| `>>> sparse.linalg.norm(I)` | Norm |

### Solving linear problems

| | |
|---|---|
| `>>> sparse.linalg.spsolve(H,I)` | 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)
```

## Matrix Functions

### Addition

| | |
|---|---|
| `>>> np.add(A,D)` | Addition |

### Subtraction

| | |
|---|---|
| `>>> np.subtract(A,D)` | Subtraction |

### Division

| | |
|---|---|
| `>>> np.divide(A,D)` | Division |

### Multiplication

| | |
|---|---|
| `>>> np.multiply(D,A)` | Multiplication |
| `>>> np.dot(A,D)` | Dot product |
| `>>> np.vdot(A,D)` | Vector dot product |
| `>>> np.inner(A,D)` | Inner product |
| `>>> np.outer(A,D)` | Outer product |
| `>>> np.tensordot(A,D)` | Tensor dot product |
| `>>> np.kron(A,D)` | Kronecker product |

### Exponential Functions

| | |
|---|---|
| `>>> linalg.expm(A)` | Matrix exponential |
| `>>> linalg.expm2(A)` | Matrix exponential (Taylor Series) |
| `>>> linalg.expm3(D)` | Matrix exponential (eigenvalue decomposition) |

### Logarithm Function

| | |
|---|---|
| `>>> linalg.logm(A)` | Matrix logarithm |

### Trigonometric Tunctions

| | |
|---|---|
| `>>> linalg.sinm(D)` | Matrix sine |
| `>>> linalg.cosm(D)` | Matrix cosine |
| `>>> linalg.tanm(A)` | Matrix tangent |

### Hyperbolic Trigonometric Functions

| | |
|---|---|
| `>>> linalg.sinhm(D)` | Hypberbolic matrix sine |
| `>>> linalg.coshm(D)` | Hyperbolic matrix cosine |
| `>>> linalg.tanhm(A)` | Hyperbolic matrix tangent |

### Matrix Sign Function

| | |
|---|---|
| `>>> np.sigm(A)` | Matrix sign function |

### Matrix Square Root

| | |
|---|---|
| `>>> linalg.sqrtm(A)` | Matrix square root |

### Arbitrary Functions

| | |
|---|---|
| `>>> linalg.funm(A, lambda x: x*x)` | Evaluate matrix function |

## Decompositions

### Eigenvalues and Eigenvectors

| | |
|---|---|
| `>>> la, v = linalg.eig(A)` | Solve ordinary or generalized eigenvalue problem for square matrix |
| `>>> l1, l2 = la` | Unpack eigenvalues |
| `>>> v[:,0]` | First eigenvector |
| `>>> v[:,1]` | Second eigenvector |
| `>>> linalg.eigvals(A)` | Unpack eigenvalues |

### Singular Value Decomposition

| | |
|---|---|
| `>>> U,s,Vh = linalg.svd(B)` | Singular Value Decomposition (SVD) |
| `>>> M,N = B.shape` | |
| `>>> Sig = linalg.diagsvd(s,M,N)` | Construct sigma matrix in SVD |

### LU Decomposition

| | |
|---|---|
| `>>> P,L,U = linalg.lu(C)` | LU Decomposition |

## Sparse Matrix Decompositions

| | |
|---|---|
| `>>> la, v = sparse.linalg.eigs(F,1)` | Eigenvalues and eigenvectors |
| `>>> sparse.linalg.svds(H, 2)` | SVD |

# Python For Data Science *Cheat Sheet*
## Seaborn

## 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")          Step 1
>>> sns.set_style("whitegrid")               Step 2
>>> g = sns.lmplot(x="tip",                  Step 3
                   y="total_bill",
                   data=tips,
                   aspect=2)
>>> g = (g.set_axis_labels("Tip","Total bill(USD)").
set(xlim=(0,10),ylim=(0,100)))
>>> plt.title("title")                       Step 4
>>> plt.show(g)                              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")
```

## 2  Figure Aesthetics
### Also see Matplotlib

| `>>> f, ax = plt.subplots(figsize=(5,6))` | Create a figure and one subplot |
|---|---|

### Seaborn styles

```
>>> sns.set()
>>> sns.set_style("whitegrid")
>>> sns.set_style("ticks",
                  {"xtick.major.size":8,
                   "ytick.major.size":8})
>>> sns.axes_style("whitegrid")
```

- (Re)set the seaborn default
- Set the matplotlib parameters
- Set the matplotlib parameters

Return a dict of params or use with `with` to temporarily set the style

### Context Functions

```
>>> sns.set_context("talk")
>>> sns.set_context("notebook",
                    font_scale=1.5,
                    rc={"lines.linewidth":2.5})
```

- Set context to `"talk"`
- Set context to `"notebook"`, scale font elements and override param mapping

### Color Palette

```
>>> sns.set_palette("husl",3)
>>> sns.color_palette("husl")
>>> flatui = ["#9b59b6","#3498db","#95a5a6","#e74c3c","#34495e","#2ecc71"]
>>> sns.set_palette(flatui)
```

- Define the color palette
- Use with `with` to temporarily set palette
- Set your own color palette

## 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

#### Scatterplot

```
>>> sns.stripplot(x="species",
                  y="petal_length",
                  data=iris)
>>> sns.swarmplot(x="species",
                  y="petal_length",
                  data=iris)
```

- Scatterplot with one categorical variable
- Categorical scatterplot with non-overlapping points

#### Bar Chart

```
>>> sns.barplot(x="sex",
                y="survived",
                hue="class",
                data=titanic)
```

- Show point estimates and confidence intervals with scatterplot glyphs

#### Count Plot

```
>>> sns.countplot(x="deck",
                  data=titanic,
                  palette="Greens_d")
```

- Show count of observations

#### Point Plot

```
>>> sns.pointplot(x="class",
                  y="survived",
                  hue="sex",
                  data=titanic,
                  palette={"male":"g",
                           "female":"m"},
                  markers=["^","o"],
                  linestyles=["-","--"])
```

- Show point estimates and confidence intervals as rectangular bars

#### Boxplot

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

- Boxplot
- Boxplot with wide-form data

#### Violinplot

```
>>> sns.violinplot(x="age",
                   y="sex",
                   hue="survived",
                   data=titanic)
```

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

- Clear an axis
- Clear an entire figure
- Close a window