# The Complete Python for Data Science Cheat Sheet Booklet

- Python Basics
- NumPy
- Pandas
- Matplotlib
- Scikit-learn
- Seaborn
- Keras
- Scipy
- PySpark

# 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
```
>>> 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
```
>>> 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]`<br>`2` | Select item at index 1 |

#### Slice
| | |
|---|---|
| `>>> my_array[0:2]`<br>`array([1, 2])` | Select items at index 0 and 1 |

#### Subset 2D Numpy arrays
| | |
|---|---|
| `>>> my_2darray[:,0]`<br>`array([1, 4])` | my_2darray[rows, columns] |

### 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*
## 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>`          [ 4. ,  5. ,  6. ]]]),`<br>`   array([[[ 3.,  2.,  3.],`<br>`          [ 4.,  5.,  6.]]])]` | Split the array vertically at the 2nd index |

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

**By Label**

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

**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)]` | Series `s` where value is not >1 |
| `>>> s[(s < -1) | (s > 2)]` | `s`  where value is <-1 or >2 |
| `>>> df[df['Population']>1200000000]` | 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()`

```
>>> df.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*
## 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.

*matplotlib*

## 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 horiontal 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
```

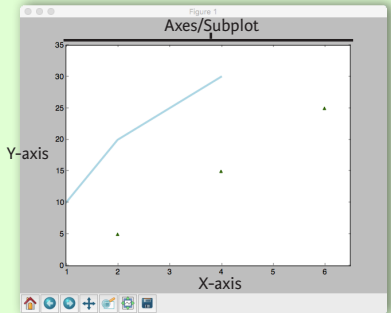### 2D Data or Images

```
>>> 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*
## Scikit-Learn

## Scikit-learn

Scikit-learn is an open source Python library that implements a range of machine learning, preprocessing, cross-validation and visualization algorithms using a unified interface.

### A Basic Example

```
>>> from sklearn import neighbors, datasets, preprocessing
>>> from sklearn.model_selection import train_test_split
>>> from sklearn.metrics import accuracy_score
>>> iris = datasets.load_iris()
>>> X, y = iris.data[:, :2], iris.target
>>> X_train, X_test, y_train, y_test = train_test_split(X,y,random_state=33)
>>> scaler = preprocessing.StandardScaler().fit(X_train)
>>> X_train = scaler.transform(X_train)
>>> X_test = scaler.transform(X_test)
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
>>> knn.fit(X_train, y_train)
>>> y_pred = knn.predict(X_test)
>>> accuracy_score(y_test, y_pred)
```

## Loading The Data                                    Also see NumPy & Pandas

Your data needs to be numeric and stored as NumPy arrays or SciPy sparse matrices. Other types that are convertible to numeric arrays, such as Pandas DataFrame, are also acceptable.

```
>>> import numpy as np
>>> X = np.random.random((10,5))
>>> y = np.array(['M','M','F','F','M','F','M','M','F','F','F'])
>>> X[X < 0.7] = 0
```

## Training And Test Data

```
>>> from sklearn.model_selection import train_test_split
>>> X_train, X_test, y_train, y_test = train_test_split(X,
                                                        y,
                                                        random_state=0)
```

## Preprocessing The Data

### Standardization

```
>>> from sklearn.preprocessing import StandardScaler
>>> scaler = StandardScaler().fit(X_train)
>>> standardized_X = scaler.transform(X_train)
>>> standardized_X_test = scaler.transform(X_test)
```

### Normalization

```
>>> from sklearn.preprocessing import Normalizer
>>> scaler = Normalizer().fit(X_train)
>>> normalized_X = scaler.transform(X_train)
>>> normalized_X_test = scaler.transform(X_test)
```

### Binarization

```
>>> from sklearn.preprocessing import Binarizer
>>> binarizer = Binarizer(threshold=0.0).fit(X)
>>> binary_X = binarizer.transform(X)
```

### Encoding Categorical Features

```
>>> from sklearn.preprocessing import LabelEncoder
>>> enc = LabelEncoder()
>>> y = enc.fit_transform(y)
```

### Imputing Missing Values

```
>>> from sklearn.preprocessing import Imputer
>>> imp = Imputer(missing_values=0, strategy='mean', axis=0)
>>> imp.fit_transform(X_train)
```

### Generating Polynomial Features

```
>>> from sklearn.preprocessing import PolynomialFeatures
>>> poly = PolynomialFeatures(5)
>>> poly.fit_transform(X)
```

## Create Your Model

### Supervised Learning Estimators

#### Linear Regression
```
>>> from sklearn.linear_model import LinearRegression
>>> lr = LinearRegression(normalize=True)
```

#### Support Vector Machines (SVM)
```
>>> from sklearn.svm import SVC
>>> svc = SVC(kernel='linear')
```

#### Naive Bayes
```
>>> from sklearn.naive_bayes import GaussianNB
>>> gnb = GaussianNB()
```

#### KNN
```
>>> from sklearn import neighbors
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
```

### Unsupervised Learning Estimators

#### Principal Component Analysis (PCA)
```
>>> from sklearn.decomposition import PCA
>>> pca = PCA(n_components=0.95)
```

#### K Means
```
>>> from sklearn.cluster import KMeans
>>> k_means = KMeans(n_clusters=3, random_state=0)
```

## Model Fitting

### Supervised learning
```
>>> lr.fit(X, y)
>>> knn.fit(X_train, y_train)
>>> svc.fit(X_train, y_train)
```
Fit the model to the data

### Unsupervised Learning
```
>>> k_means.fit(X_train)
>>> pca_model = pca.fit_transform(X_train)
```
Fit the model to the data
Fit to data, then transform it

## Prediction

### Supervised Estimators
```
>>> y_pred = svc.predict(np.random.random((2,5)))
>>> y_pred = lr.predict(X_test)
>>> y_pred = knn.predict_proba(X_test)
```
Predict labels
Predict labels
Estimate probability of a label

### Unsupervised Estimators
```
>>> y_pred = k_means.predict(X_test)
```
Predict labels in clustering algos

## Evaluate Your Model's Performance

### Classification Metrics

#### Accuracy Score
```
>>> knn.score(X_test, y_test)
>>> from sklearn.metrics import accuracy_score
>>> accuracy_score(y_test, y_pred)
```
Estimator score method
Metric scoring functions

#### Classification Report
```
>>> from sklearn.metrics import classification_report
>>> print(classification_report(y_test, y_pred))
```
Precision, recall, f1-score and support

#### Confusion Matrix
```
>>> from sklearn.metrics import confusion_matrix
>>> print(confusion_matrix(y_test, y_pred))
```

### Regression Metrics

#### Mean Absolute Error
```
>>> from sklearn.metrics import mean_absolute_error
>>> y_true = [3, -0.5, 2]
>>> mean_absolute_error(y_true, y_pred)
```

#### Mean Squared Error
```
>>> from sklearn.metrics import mean_squared_error
>>> mean_squared_error(y_test, y_pred)
```

#### $R^2$ Score
```
>>> from sklearn.metrics import r2_score
>>> r2_score(y_true, y_pred)
```

### Clustering Metrics

#### Adjusted Rand Index
```
>>> from sklearn.metrics import adjusted_rand_score
>>> adjusted_rand_score(y_true, y_pred)
```

#### Homogeneity
```
>>> from sklearn.metrics import homogeneity_score
>>> homogeneity_score(y_true, y_pred)
```

#### V-measure
```
>>> from sklearn.metrics import v_measure_score
>>> metrics.v_measure_score(y_true, y_pred)
```

### Cross-Validation
```
>>> from sklearn.cross_validation import cross_val_score
>>> print(cross_val_score(knn, X_train, y_train, cv=4))
>>> print(cross_val_score(lr, X, y, cv=2))
```

## Tune Your Model

### Grid Search
```
>>> from sklearn.grid_search import GridSearchCV
>>> params = {"n_neighbors": np.arange(1,3),
              "metric": ["euclidean", "cityblock"]}
>>> grid = GridSearchCV(estimator=knn,
                        param_grid=params)
>>> grid.fit(X_train, y_train)
>>> print(grid.best_score_)
>>> print(grid.best_estimator_.n_neighbors)
```

### Randomized Parameter Optimization
```
>>> from sklearn.grid_search import RandomizedSearchCV
>>> params = {"n_neighbors": range(1,5),
              "weights": ["uniform", "distance"]}
>>> rsearch = RandomizedSearchCV(estimator=knn,
                                 param_distributions=params,
                                 cv=4,
                                 n_iter=8,
                                 random_state=5)
>>> rsearch.fit(X_train, y_train)
>>> print(rsearch.best_score_)
```

# 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)))               Step 4
>>> plt.title("title")
>>> 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()                              (Re)set the seaborn default
>>> sns.set_style("whitegrid")             Set the matplotlib parameters
>>> sns.set_style("ticks",                 Set the matplotlib parameters
                  {"xtick.major.size":8,
                   "ytick.major.size":8})
>>> sns.axes_style("whitegrid")            Return a dict of params or use with
                                           with to temporarily set the style
```

### Context Functions

```
>>> sns.set_context("talk")                Set context to "talk"
>>> sns.set_context("notebook",            Set context to "notebook",
                    font_scale=1.5,        scale font elements and
                    rc={"lines.linewidth":2.5})  override param mapping
```

### Color Palette

```
>>> sns.set_palette("husl",3)              Define the color palette
>>> sns.color_palette("husl")              Use with with to temporarily set palette
>>> flatui = ["#9b59b6","#3498db","#95a5a6","#e74c3c","#34495e","#2ecc71"]
>>> sns.set_palette(flatui)                Set your own color palette
```

## 3 Plotting With Seaborn

### Axis Grids

```
>>> g = sns.FacetGrid(titanic,
                      col="survived",      Subplot grid for plotting conditional
                      row="sex")           relationships
>>> g = g.map(plt.hist,"age")
>>> sns.factorplot(x="pclass",             Draw a categorical plot onto a
                   y="survived",           Facetgrid
                   hue="sex",
                   data=titanic)
>>> sns.lmplot(x="sepal_width",            Plot data and regression model fits
               y="sepal_length",          across a FacetGrid
               hue="species",
               data=iris)
```

```
>>> h = sns.PairGrid(iris)                 Subplot grid for plotting pairwise
>>> h = h.map(plt.scatter)                 relationships
>>> sns.pairplot(iris)                     Plot pairwise bivariate distributions
>>> i = sns.JointGrid(x="x",               Grid for bivariate plot with marginal
                      y="y",               univariate plots
                      data=data)
>>> i = i.plot(sns.regplot,
               sns.distplot)
>>> sns.jointplot("sepal_length",          Plot bivariate distribution
                  "sepal_width",
                  data=iris,
                  kind='kde')
```

### Categorical Plots

#### Scatterplot

```
>>> sns.stripplot(x="species",             Scatterplot with one
                  y="petal_length",        categorical variable
                  data=iris)
>>> sns.swarmplot(x="species",             Categorical scatterplot with
                  y="petal_length",        non-overlapping points
                  data=iris)
```

#### Bar Chart

```
>>> sns.barplot(x="sex",                   Show point estimates and
                y="survived",              confidence intervals with
                hue="class",               scatterplot glyphs
                data=titanic)
```

#### Count Plot

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

#### Point Plot

```
>>> sns.pointplot(x="class",               Show point estimates and
                  y="survived",            confidence intervals as
                  hue="sex",               rectangular bars
                  data=titanic,
                  palette={"male":"g",
                           "female":"m"},
                  markers=["^","o"],
                  linestyles=["-","--"])
```

#### Boxplot

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

#### Violinplot

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

### Regression Plots

```
>>> sns.regplot(x="sepal_width",           Plot data and a linear regression
                y="sepal_length",          model fit
                data=iris,
                ax=ax)
```

### Distribution Plots

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

### Matrix Plots

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

## 4 Further Customizations     Also see Matplotlib

### Axisgrid Objects

```
>>> g.despine(left=True)                   Remove left spine
>>> g.set_ylabels("Survived")              Set the labels of the y-axis
>>> g.set_xticklabels(rotation=45)         Set the tick labels for x
>>> g.set_axis_labels("Survived",          Set the axis labels
                      "Sex")
>>> h.set(xlim=(0,5),                      Set the limit and ticks of the
          ylim=(0,5),                      x-and y-axis
          xticks=[0,2.5,5],
          yticks=[0,2.5,5])
```

### Plot

```
>>> plt.title("A Title")                   Add plot title
>>> plt.ylabel("Survived")                 Adjust the label of the y-axis
>>> plt.xlabel("Sex")                      Adjust the label of the x-axis
>>> plt.ylim(0,100)                        Adjust the limits of the y-axis
>>> plt.xlim(0,10)                         Adjust the limits of the x-axis
>>> plt.setp(ax,yticks=[0,5])              Adjust a plot property
>>> plt.tight_layout()                     Adjust subplot params
```

## 5 Show or Save Plot     Also see Matplotlib

```
>>> plt.show()                             Show the plot
>>> plt.savefig("foo.png")                 Save the plot as a figure
>>> plt.savefig("foo.png",                 Save transparent figure
                transparent=True)
```

### Close & Clear     Also see Matplotlib

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

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

## Keras

**Keras** is a powerful and easy-to-use deep learning library for Theano and TensorFlow that provides a high-level neural networks API to develop and evaluate deep learning models.

### A Basic Example

```python
>>> import numpy as np
>>> from keras.models import Sequential
>>> from keras.layers import Dense
>>> data = np.random.random((1000,100))
>>> labels = np.random.randint(2,size=(1000,1))
>>> model = Sequential()
>>> model.add(Dense(32,
                    activation='relu',
                    input_dim=100))
>>> model.add(Dense(1, activation='sigmoid'))
>>> model.compile(optimizer='rmsprop',
                  loss='binary_crossentropy',
                  metrics=['accuracy'])
>>> model.fit(data,labels,epochs=10,batch_size=32)
>>> predictions4 = model.predict(data)
```

### Data                     Also see NumPy, Pandas & Scikit-Learn

Your data needs to be stored as NumPy arrays or as a list of NumPy arrays. Ideally, you split the data in training and test sets, for which you can also resort to the `train_test_split` module of `sklearn.cross_validation`.

### Keras Data Sets

```python
>>> from keras.datasets import boston_housing,
                                mnist,
                                cifar10,
                                imdb
>>> (x_train,y_train),(x_test,y_test) = mnist.load_data()
>>> (x_train2,y_train2),(x_test2,y_test2) = boston_housing.load_data()
>>> (x_train3,y_train3),(x_test3,y_test3) = cifar10.load_data()
>>> (x_train4,y_train4),(x_test4,y_test4) = imdb.load_data(num_words=20000)
>>> num_classes = 10
```

### Other

```python
>>> from urllib.request import urlopen
>>> data = np.loadtxt(urlopen("http://archive.ics.uci.edu/
ml/machine-learning-databases/pima-indians-diabetes/
pima-indians-diabetes.data"),delimiter=",")
>>> X = data[:,0:8]
>>> y = data [:,8]
```

### Preprocessing                     Also see NumPy & Scikit-Learn

### Sequence Padding

```python
>>> from keras.preprocessing import sequence
>>> x_train4 = sequence.pad_sequences(x_train4,maxlen=80)
>>> x_test4 = sequence.pad_sequences(x_test4,maxlen=80)
```

### One-Hot Encoding

```python
>>> from keras.utils import to_categorical
>>> Y_train = to_categorical(y_train, num_classes)
>>> Y_test = to_categorical(y_test, num_classes)
>>> Y_train3 = to_categorical(y_train3, num_classes)
>>> Y_test3 = to_categorical(y_test3, num_classes)
```

### Train and Test Sets

```python
>>> from sklearn.model_selection import train_test_split
>>> X_train5,X_test5,y_train5,y_test5 = train_test_split(X,
                                                         y,
                                                         test_size=0.33,
                                                         random_state=42)
```

### Standardization/Normalization

```python
>>> from sklearn.preprocessing import StandardScaler
>>> scaler = StandardScaler().fit(x_train2)
>>> standardized_X = scaler.transform(x_train2)
>>> standardized_X_test = scaler.transform(x_test2)
```

## Model Architecture

### Sequential Model

```python
>>> from keras.models import Sequential
>>> model = Sequential()
>>> model2 = Sequential()
>>> model3 = Sequential()
```

### Multilayer Perceptron (MLP)

#### Binary Classification

```python
>>> from keras.layers import Dense
>>> model.add(Dense(12,
                    input_dim=8,
                    kernel_initializer='uniform',
                    activation='relu'))
>>> model.add(Dense(8,kernel_initializer='uniform',activation='relu'))
>>> model.add(Dense(1,kernel_initializer='uniform',activation='sigmoid'))
```

#### Multi-Class Classification

```python
>>> from keras.layers import Dropout
>>> model.add(Dense(512,activation='relu',input_shape=(784,)))
>>> model.add(Dropout(0.2))
>>> model.add(Dense(512,activation='relu'))
>>> model.add(Dropout(0.2))
>>> model.add(Dense(10,activation='softmax'))
```

#### Regression

```python
>>> model.add(Dense(64,activation='relu',input_dim=train_data.shape[1]))
>>> model.add(Dense(1))
```

### Convolutional Neural Network (CNN)

```python
>>> from keras.layers import Activation,Conv2D,MaxPooling2D,Flatten
>>> model2.add(Conv2D(32,(3,3),padding='same',input_shape=x_train.shape[1:]))
>>> model2.add(Activation('relu'))
>>> model2.add(Conv2D(32,(3,3)))
>>> model2.add(Activation('relu'))
>>> model2.add(MaxPooling2D(pool_size=(2,2)))
>>> model2.add(Dropout(0.25))

>>> model2.add(Conv2D(64,(3,3), padding='same'))
>>> model2.add(Activation('relu'))
>>> model2.add(Conv2D(64,(3, 3)))
>>> model2.add(Activation('relu'))
>>> model2.add(MaxPooling2D(pool_size=(2,2)))
>>> model2.add(Dropout(0.25))

>>> model2.add(Flatten())
>>> model2.add(Dense(512))
>>> model2.add(Activation('relu'))
>>> model2.add(Dropout(0.5))
>>> model2.add(Dense(num_classes))
>>> model2.add(Activation('softmax'))
```

### Recurrent Neural Network (RNN)

```python
>>> from keras.klayers import Embedding,LSTM
>>> model3.add(Embedding(20000,128))
>>> model3.add(LSTM(128,dropout=0.2,recurrent_dropout=0.2))
>>> model3.add(Dense(1,activation='sigmoid'))
```

## Inspect Model

```python
>>> model.output_shape        Model output shape
>>> model.summary()           Model summary representation
>>> model.get_config()        Model configuration
>>> model.get_weights()       List all weight tensors in the model
```

## Compile Model

#### MLP: Binary Classification

```python
>>> model.compile(optimizer='adam',
                  loss='binary_crossentropy',
                  metrics=['accuracy'])
```

#### MLP: Multi-Class Classification

```python
>>> model.compile(optimizer='rmsprop',
                  loss='categorical_crossentropy',
                  metrics=['accuracy'])
```

#### MLP: Regression

```python
>>> model.compile(optimizer='rmsprop',
                  loss='mse',
                  metrics=['mae'])
```

#### Recurrent Neural Network

```python
>>> model3.compile(loss='binary_crossentropy',
                   optimizer='adam',
                   metrics=['accuracy'])
```

## Model Training

```python
>>> model3.fit(x_train4,
               y_train4,
               batch_size=32,
               epochs=15,
               verbose=1,
               validation_data=(x_test4,y_test4))
```

## Evaluate Your Model's Performance

```python
>>> score = model3.evaluate(x_test,
                            y_test,
                            batch_size=32)
```

## Prediction

```python
>>> model3.predict(x_test4, batch_size=32)
>>> model3.predict_classes(x_test4,batch_size=32)
```

## Save/ Reload Models

```python
>>> from keras.models import load_model
>>> model3.save('model_file.h5')
>>> my_model = load_model('my_model.h5')
```

## Model Fine-tuning

### Optimization Parameters

```python
>>> from keras.optimizers import RMSprop
>>> opt = RMSprop(lr=0.0001, decay=1e-6)
>>> model2.compile(loss='categorical_crossentropy',
                   optimizer=opt,
                   metrics=['accuracy'])
```

### Early Stopping

```python
>>> from keras.callbacks import EarlyStopping
>>> early_stopping_monitor = EarlyStopping(patience=2)
>>> model3.fit(x_train4,
               y_train4,
               batch_size=32,
               epochs=15,
               validation_data=(x_test4,y_test4),
               callbacks=[early_stopping_monitor])
```

# 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

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

### Shape Manipulation

```
>>> np.transpose(b)
>>> b.flatten()
>>> np.hstack((b,c))
>>> np.vstack((a,b))
>>> np.hsplit(c,2)
>>> np.vpslit(d,2)
```
Permute array dimensions
Flatten the array
Stack arrays horizontally (column-wise)
Stack arrays vertically (row-wise)
Split the array horizontally at the 2nd index
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)
>>> np.imag(c)
>>> np.real_if_close(c,tol=1000)
>>> np.cast['f'](np.pi)
```
Return the real part of the array elements
Return the imaginary part of the array elements
Return a real array if complex parts close to 0
Cast object to a data type

### Other Useful Functions

```
>>> np.angle(b,deg=True)
>>> g = np.linspace(0,np.pi,num=5)

>>> g [3:] += np.pi
>>> np.unwrap(g)
>>> np.logspace(0,10,3)
>>> np.select([c<4],[c*2])

>>> misc.factorial(a)
>>> misc.comb(10,3,exact=True)
>>> misc.central_diff_weights(3)
>>> misc.derivative(myfunc,1.0)
```
Return the angle of the complex argument
Create an array of evenly spaced values (number of samples)
Unwrap
Create an array of evenly spaced values (log scale)
Return values from a list of arrays depending on conditions
Factorial
Combine N things taken at k time
Weights for Np-point central derivative
Find the n-th derivative of a function at a point

## 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                    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*
## PySpark - RDD Basics

## Spark

**PySpark** is the Spark Python API that exposes the Spark programming model to Python.

## Initializing Spark

### SparkContext

```
>>> from pyspark import SparkContext
>>> sc = SparkContext(master = 'local[2]')
```

### Inspect SparkContext

| | |
|---|---|
| `>>> sc.version` | Retrieve SparkContext version |
| `>>> sc.pythonVer` | Retrieve Python version |
| `>>> sc.master` | Master URL to connect to |
| `>>> str(sc.sparkHome)` | Path where Spark is installed on worker nodes |
| `>>> str(sc.sparkUser())` | Retrieve name of the Spark User running SparkContext |
| `>>> sc.appName` | Return application name |
| `>>> sc.applicationId` | Retrieve application ID |
| `>>> sc.defaultParallelism` | Return default level of parallelism |
| `>>> sc.defaultMinPartitions` | Default minimum number of partitions for RDDs |

### Configuration

```
>>> from pyspark import SparkConf, SparkContext
>>> conf = (SparkConf()
             .setMaster("local")
             .setAppName("My app")
             .set("spark.executor.memory", "1g"))
>>> sc = SparkContext(conf = conf)
```

### Using The Shell

In the PySpark shell, a special interpreter-aware SparkContext is already created in the variable called `sc`.

```
$ ./bin/spark-shell --master local[2]
$ ./bin/pyspark --master local[4] --py-files code.py
```

Set which master the context connects to with the `--master` argument, and add Python .zip, .egg or .py files to the runtime path by passing a comma-separated list to `--py-files`.

## Loading Data

### Parallelized Collections

```
>>> rdd = sc.parallelize([('a',7),('a',2),('b',2)])
>>> rdd2 = sc.parallelize([('a',2),('d',1),('b',1)])
>>> rdd3 = sc.parallelize(range(100))
>>> rdd4 = sc.parallelize([("a",["x","y","z"]),
                           ("b",["p", "r"])])
```

### External Data

Read either one text file from HDFS, a local file system or or any Hadoop-supported file system URI with `textFile()`, or read in a directory of text files with `wholeTextFiles()`.

```
>>> textFile = sc.textFile("/my/directory/*.txt")
>>> textFile2 = sc.wholeTextFiles("/my/directory/")
```

## Retrieving RDD Information

### Basic Information

| | |
|---|---|
| `>>> rdd.getNumPartitions()` | List the number of partitions |
| `>>> rdd.count()`<br>`3` | Count RDD instances |
| `>>> rdd.countByKey()`<br>`defaultdict(<type 'int'>,{'a':2,'b':1})` | Count RDD instances by key |
| `>>> rdd.countByValue()`<br>`defaultdict(<type 'int'>,{('b',2):1,('a',2):1,('a',7):1})` | Count RDD instances by value |
| `>>> rdd.collectAsMap()`<br>`{'a': 2,'b': 2}` | Return (key,value) pairs as a dictionary |
| `>>> rdd3.sum()`<br>`4950` | Sum of RDD elements |
| `>>> sc.parallelize([]).isEmpty()`<br>`True` | Check whether RDD is empty |

### Summary

| | |
|---|---|
| `>>> rdd3.max()`<br>`99` | Maximum value of RDD elements |
| `>>> rdd3.min()`<br>`0` | Minimum value of RDD elements |
| `>>> rdd3.mean()`<br>`49.5` | Mean value of RDD elements |
| `>>> rdd3.stdev()`<br>`28.866070047722118` | Standard deviation of RDD elements |
| `>>> rdd3.variance()`<br>`833.25` | Compute variance of RDD elements |
| `>>> rdd3.histogram(3)`<br>`([0,33,66,99],[33,33,34])` | Compute histogram by bins |
| `>>> rdd3.stats()` | Summary statistics (count, mean, stdev, max & min) |

## Applying Functions

| | |
|---|---|
| `>>> rdd.map(lambda x: x+(x[1],x[0]))`<br>`    .collect()`<br>`[('a',7,7,'a'),('a',2,2,'a'),('b',2,2,'b')]` | Apply a function to each RDD element |
| `>>> rdd5 = rdd.flatMap(lambda x: x+(x[1],x[0]))`<br>`>>> rdd5.collect()`<br>`['a',7,7,'a','a',2,2,'a','b',2,2,'b']` | Apply a function to each RDD element and flatten the result |
| `>>> rdd4.flatMapValues(lambda x: x)`<br>`    .collect()`<br>`[('a','x'),('a','y'),('a','z'),('b','p'),('b','r')]` | Apply a flatMap function to each (key,value) pair of `rdd4` without changing the keys |

## Selecting Data

### Getting

| | |
|---|---|
| `>>> rdd.collect()`<br>`[('a', 7), ('a', 2), ('b', 2)]` | Return a list with all RDD elements |
| `>>> rdd.take(2)`<br>`[('a', 7), ('a', 2)]` | Take first 2 RDD elements |
| `>>> rdd.first()`<br>`('a', 7)` | Take first RDD element |
| `>>> rdd.top(2)`<br>`[('b', 2), ('a', 7)]` | Take top 2 RDD elements |

### Sampling

| | |
|---|---|
| `>>> rdd3.sample(False, 0.15, 81).collect()`<br>`[3,4,27,31,40,41,42,43,60,76,79,80,86,97]` | Return sampled subset of `rdd3` |

### Filtering

| | |
|---|---|
| `>>> rdd.filter(lambda x: "a" in x)`<br>`    .collect()`<br>`[('a',7),('a',2)]` | Filter the RDD |
| `>>> rdd5.distinct().collect()`<br>`['a',2,'b',7]` | Return distinct RDD values |
| `>>> rdd.keys().collect()`<br>`['a', 'a', 'b']` | Return (key,value) RDD's keys |

## Iterating

| | |
|---|---|
| `>>> def g(x): print(x)`<br>`>>> rdd.foreach(g)`<br>`('a', 7)`<br>`('b', 2)`<br>`('a', 2)` | Apply a function to all RDD elements |

## Reshaping Data

### Reducing

| | |
|---|---|
| `>>> rdd.reduceByKey(lambda x,y : x+y)`<br>`    .collect()`<br>`[('a',9),('b',2)]` | Merge the `rdd` values for each key |
| `>>> rdd.reduce(lambda a, b: a + b)`<br>`('a',7,'a',2,'b',2)` | Merge the `rdd` values |

### Grouping by

| | |
|---|---|
| `>>> rdd3.groupBy(lambda x: x % 2)`<br>`    .mapValues(list)`<br>`    .collect()` | Return RDD of grouped values |
| `>>> rdd.groupByKey()`<br>`    .mapValues(list)`<br>`    .collect()`<br>`[('a',[7,2]),('b',[2])]` | Group `rdd` by key |

### Aggregating

| | |
|---|---|
| `>>> seqOp = (lambda x,y: (x[0]+y,x[1]+1))` | |
| `>>> combOp = (lambda x,y:(x[0]+y[0],x[1]+y[1]))` | |
| `>>> rdd3.aggregate((0,0),seqOp,combOp)`<br>`(4950,100)` | Aggregate RDD elements of each partition and then the results |
| `>>> rdd.aggregateByKey((0,0),seqop,combop)`<br>`    .collect()`<br>`[('a',(9,2)), ('b',(2,1))]` | Aggregate values of each RDD key |
| `>>> rdd3.fold(0,add)`<br>`4950` | Aggregate the elements of each partition, and then the results |
| `>>> rdd.foldByKey(0, add)`<br>`    .collect()`<br>`[('a',9), ('b',2)]` | Merge the values for each key |
| `>>> rdd3.keyBy(lambda x: x+x)`<br>`    .collect()` | Create tuples of RDD elements by applying a function |

## Mathematical Operations

| | |
|---|---|
| `>>> rdd.subtract(rdd2)`<br>`    .collect()`<br>`[('b',2),('a',7)]` | Return each `rdd` value not contained in `rdd2` |
| `>>> rdd2.subtractByKey(rdd)`<br>`    .collect()`<br>`[('d', 1)]` | Return each (key,value) pair of `rdd2` with no matching key in `rdd` |
| `>>> rdd.cartesian(rdd2).collect()` | Return the Cartesian product of `rdd` and `rdd2` |

## Sort

| | |
|---|---|
| `>>> rdd2.sortBy(lambda x: x[1])`<br>`    .collect()`<br>`[('d',1),('b',1),('a',2)]` | Sort RDD by given function |
| `>>> rdd2.sortByKey()`<br>`    .collect()`<br>`[('a',2),('b',1),('d',1)]` | Sort (key, value) RDD by key |

## Repartitioning

| | |
|---|---|
| `>>> rdd.repartition(4)` | New RDD with 4 partitions |
| `>>> rdd.coalesce(1)` | Decrease the number of partitions in the RDD to 1 |

## Saving

```
>>> rdd.saveAsTextFile("rdd.txt")
>>> rdd.saveAsHadoopFile("hdfs://namenodehost/parent/child",
          'org.apache.hadoop.mapred.TextOutputFormat')
```

## Stopping SparkContext

```
>>> sc.stop()
```

## Execution

```
$ ./bin/spark-submit examples/src/main/python/pi.py
```