

CHEAT SHEET'S

DATA SCIENCE, MACHINE LEARNING & NEURAL NETWORKS.

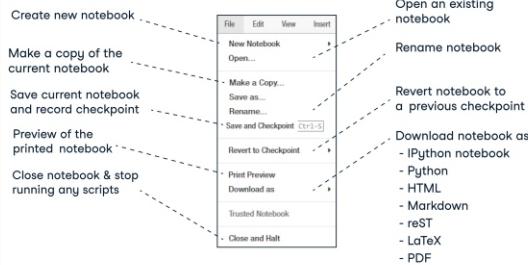
- | | |
|--|--|
| <ol style="list-style-type: none">1. Jupyter Notebook2. Python (high-level, general-purpose programming language)3. Sql (Structured Query Language)4. Numpy (N-dimentional arrays)5. Pandas (Series, DataFrames)6. Matplotlib (Visualization)7. Seaborn (Visualization)8. Scipy (Visualization)9. Bokeh (Visualization)10. Folium (Map-Visualization) | <ol style="list-style-type: none">11. Scikit-Learn (Model Building)12. Keras (Creating deep models)13. spaCy (pre-process text for deep learning)14. Pyspark SQL (API)15. Pyspark RDD (Parallel processing on a cluster)16. BeautifulSoup (Web Scraping)17. Selenium (Web Scraping)18. Data Wrangling19. Importing Data20. Project Steps (Data Extraction, Data Cleaning, Data Wrangling, Analysis, Action) |
|--|--|

~ By G PREMSAGAR

Python For Data Science

Jupyter Cheat Sheet

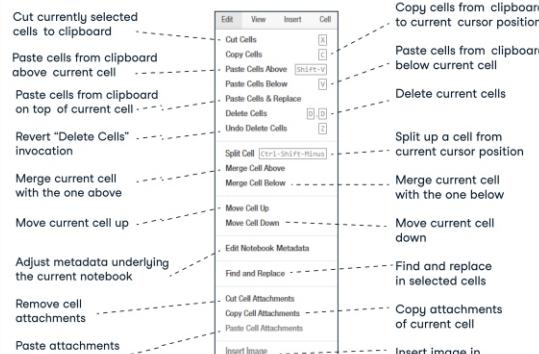
> Saving/Loading Notebooks



> Writing Code And Text

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

Edit Cells



Insert Cells



> Working with Different Programming Languages

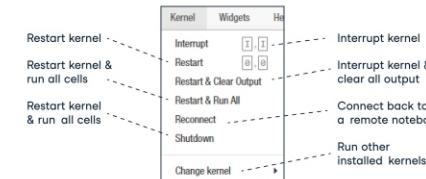
Kernels provide computation and communication with front-end interfaces like the notebooks.
There are three main kernels:

IP[y]:
IPython

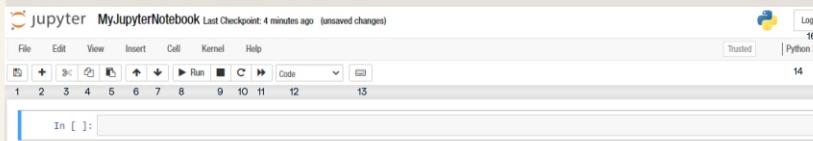
IR:
IRkernel

IJulia:
IJulia

Installing Jupyter Notebook will automatically install the IPython kernel.



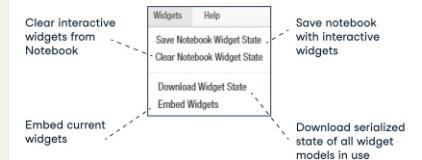
Command Mode:



> Widgets

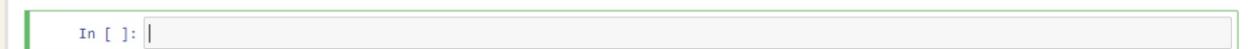
Notebook widgets provide the ability to visualize and control changes in your data, often as a control like a slider, textbox, etc.

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

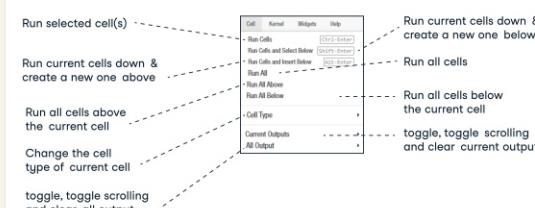


1. Save and checkpoint
2. Insert cell below
3. Cut cell
4. Copy cell(s)
5. Paste cell(s) below
6. Move cell up
7. Move cell down
8. Run current cell
9. Interrupt kernel
10. Restart kernel
11. Restart kernel and re-run notebook
12. Display characteristics
13. Open command palette
14. Current kernel
15.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

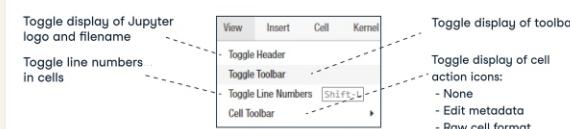
Edit Mode:



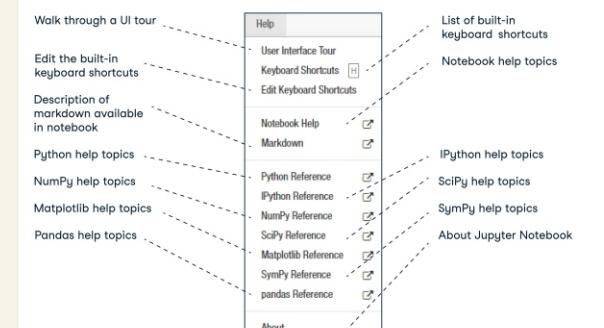
Executing Cells



View Cells



Asking For Help



Python For Data Science

python™ Basics Cheat Sheet

Variables and Data Types

Variable Assignment

```
>>> x=5
>>> x
5
```

Calculations With Variables

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

Types and Type Conversion

<code>str() '5', '3.45', 'True' #Variables to strings</code>	<code>float() 5.0, 1.0 #Variables to floats</code>
<code>int() 5, 3, 1 #Variables to integers</code>	<code>bool() True, True, True #Variables to booleans</code>

Libraries



Import Libraries

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

Selective import

```
>>> from math import pi
```

Strings

```
>>> my_string = 'thisStringIsAwesome'
>>> my_string
'thisStringIsAwesome'
```

String Operations

```
>>> my_string * 2
'thisStringIsAwesomethisStringIsAwesome'
>>> my_string + 'Init'
'thisStringIsAwesomeInit'
>>> 'i' in my_string
True
```

String Indexing

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
```

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
```

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
```

Subset Lists of Lists

```
>>> my_list2[1][0] #my_list[list][itemOfList]
>>> my_list2[1][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('i') #Append an item at a time
>>> my_list.remove('i') #Remove an item
>>> del(my_list[0:1]) #Remove an item
>>> my_list.reverse() #Reverse the list
>>> my_list.extend('i') #Append an item
>>> my_list.pop(-1) #Remove an item
>>> my_list.insert(0,'i') #Insert an item
>>> my_list.sort() #Sort the list
```

Python IDEs (Integrated Development Environment)

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

Asking For Help

```
>>> help(str)
```

~SagaR

SQL COMMANDS CHEAT SHEET

SQL Commands

The commands in SQL are called Queries and they are of two types:

- Data Definition Query:** The statements which defines the structure of a database, create tables, specify their keys, indexes and so on
- Data manipulation queries:** These are the queries which can be edited.
E.g.: Select, update and insert operation

SQL Commands			Commands and syntax for querying data from Single Table			Commands and syntax for querying data from Multiple Table		
Command	Syntax	Description	Command	Syntax	Description	Command	Syntax	Description
GROUP BY	SELECT column_name, COUNT(*) FROM table_name GROUP BY column_name;	It is an clause in SQL used for aggregate functions in collaboration with the SELECT statement	OUTER JOIN	SELECT column_name(s) FROM table_1 LEFT JOIN table_2 ON table_1.column_name = table_2.column_name;	It is sued to combine rows from different tables even if the condition is NOT TRUE	SELECT c1 FROM t	SELECT c1, c2 FROM t1 INNER JOIN t2 on condition	Select column c1 and c2 from table t1 and perform an inner join between t1 and t2
HAVING	SELECT column_name, COUNT(*) FROM table_name GROUP BY column_name HAVING COUNT(*) > value;	It is used in SQL because the WHERE keyword cannot be used in aggregating functions	ROUND	SELECT ROUND(column_name, integer) FROM table_name;	It is a function that takes the column name and a integer as an argument, and rounds the values in a column to the number of decimal places specified by an integer	SELECT * FROM t	SELECT c1, c2 FROM t1 LEFT JOIN t2 on condition	Select column c1 and c2 from table t1 and perform a left join between t1 and t2
INNER JOIN	SELECT column_name(s) FROM table_1 JOIN table_2 ON table_1.column_name = table_2.column_name;	It is used to combine rows from different tables if the Join condition goes TRUE	SELECT	SELECT column_name FROM table_name;	It is a statement that is used to fetch data from a database	SELECT c1 FROM t WHERE c1 = 'test'	SELECT c1, c2 FROM t1 RIGHT JOIN t2 on condition	Select column c1 and c2 from table t1 and perform a right join between t1 and t2
INSERT	INSERT INTO table_name (column_1, column_2, column_3) VALUES (value_1, value_2, value_3);	It is used to add new rows to a table	SELECT DISTINCT	SELECT DISTINCT column_name FROM table_name;	It is used to specify that the statement is a query which returns unique values in specified columns	SELECT c1 FROM t ORDER BY c1 ASC (DESC)	SELECT c1, c2 FROM t1 FULL OUTER JOIN t2 on condition	Select column c1 and c2 from table t1 and perform a full outer join between t1 and t2
AND	SELECT column_name(s) FROM table_name WHERE column_1 = value_1 AND column_2 = value_2;	It is an operator that is used to combine two conditions	LIMIT	SELECT column_name(s) FROM table_name WHERE column_name LIKE pattern;	It is an special operator used with the WHERE clause to search for a specific pattern in a column	UPDATE	UPDATE table_name SET some_column = some_value WHERE some_column = some_value;	It is used to edit rows in a table
AS	SELECT column_name AS 'Alias' FROM table_name;	It is an keyword in SQL that is used to rename a column or table using an alias name	MAX	SELECT MAX(column_name) FROM table_name;	It is a function that takes number of columns as an argument and return the largest value among them	WHERE	SELECT column_name(s) FROM table_name WHERE column_name operator value;	It is a clause used to filter the result set to include the rows which where the condition is TRUE
BETWEEN	SELECT column_name(s) FROM table_name WHERE column_name BETWEEN value_1 AND value_2;	It is an operator used to filter the result within a certain range	MIN	SELECT MIN(column_name) FROM table_name;	It is a function that takes number of columns as an argument and return the smallest value among them	WITH	WITH temporary_name AS (SELECT * FROM table_name) SELECT * FROM temporary_name WHERE column_name operator value;	It is used to store the result of a particular query in a temporary table using an alias
CASE	SELECT column_name, CASE WHEN condition THEN 'Result_1' WHEN condition THEN 'Result_2' ELSE 'Result_3' END FROM table_name;	It is a statement used to create different outputs inside a SELECT statement	OR	SELECT column_name FROM table_name WHERE column_name = value_1 OR column_name = value_2;	It is an operator that is used to filter the result set to contain only the rows where either condition is TRUE	DELETE	DELETE FROM table_name WHERE some_column = some_value;	It is used to remove the rows from a table
COUNT	SELECT COUNT(column_name) FROM table_name;	It is a function that takes the name of a column as argument and counts the number of rows when the column is not NULL	ORDER BY	SELECT column_name FROM table_name ORDER BY column_name ASC DESC;	It is a clause used to sort the result set by a particular column either numerically or alphabetically	AVG	SELECT AVG(column_name) FROM table_name;	It is used to aggregate a numeric column and return its average
Create TABLE	CREATE TABLE table_name (column_1 datatype, column_2 datatype, column_3 datatype);	It is used to create a new table in a database and specify the name of the table and columns inside it						

~SagaR

SQL JOINS Cheat Sheet

~SagaR

COLUMN AND TABLE ALIASES

Aliases give a temporary name to a **table** or a **column** in a table.

CAT AS c				OWNER AS o	
cat_id	cat_name	mom_id	owner_id	id	name
1	Kitty	5	1	1	John Smith
2	Hugo	1	2	2	Danielle Davis
3	Sam	2	2		
4	Misty	1	NULL		

A **column alias** renames a column in the result. A **table alias** renames a table within the query. If you define a table alias, you must use it instead of the table name everywhere in the query. The **AS** keyword is optional in defining aliases.

```
SELECT
    o.name AS owner_name,
    c.cat_name
FROM cat AS c
JOIN owner AS o
    ON c.owner_id = o.id;
```

cat_name	owner_name
Kitty	John Smith
Sam	Danielle Davis
Hugo	Danielle Davis

SELF JOIN

You can join a table to itself, for example, to show a parent-child relationship.

CAT AS child				CAT AS mom			
cat_id	cat_name	owner_id	mom_id	cat_id	cat_name	owner_id	mom_id
1	Kitty	1	5	1	Kitty	1	5
2	Hugo	2	1	2	Hugo	2	1
3	Sam	2	2	3	Sam	2	2
4	Misty	NULL	1	4	Misty	NULL	1

Each occurrence of the table must be given a **different alias**. Each column reference must be preceded with an **appropriate table alias**.

```
SELECT
    child.cat_name AS child_name,
    mom.cat_name AS mom_name
FROM cat AS child
JOIN cat AS mom
    ON child.mom_id = mom.cat_id;
```

child_name	mom_name
Hugo	Kitty
Sam	Hugo
Misty	Kitty

NON-EQUI SELF JOIN

You can use a **non-equality** in the **ON** condition, for example, to show **all different pairs** of rows.

TOY AS a			TOY AS b		
toy_id	toy_name	cat_id	cat_id	toy_id	toy_name
3	mouse	1	1	3	mouse
5	ball	1	1	5	ball
1	ball	3	3	1	ball
4	mouse	4	4	4	mouse
2	spring	NULL	NULL	2	spring

```
SELECT
    a.toy_name AS toy_a,
    b.toy_name AS toy_b
FROM toy a
JOIN toy b
    ON a.cat_id < b.cat_id;
```

cat_a_id	toy_a	cat_b_id	toy_b
1	mouse	3	ball
1	ball	3	ball
1	mouse	4	mouse
1	ball	4	mouse
3	ball	4	mouse

MULTIPLE JOINS

You can join more than two tables together. First, two tables are joined, then the third table is joined to the result of the previous joining.

TOY AS t		
toy_id	toy_name	cat_id
1	ball	3
2	spring	NULL
3	mouse	1
4	mouse	4
5	ball	1

CAT AS c				OWNER AS o	
cat_id	cat_name	mom_id	owner_id	id	name
1	Kitty	5	1	1	John Smith
2	Hugo	1	2	2	Danielle Davis
3	Sam	2	2	2	Danielle Davis
4	Misty	1	NULL		

JOIN & JOIN

```
SELECT
    t.toy_name,
    c.cat_name,
    o.name AS owner_name
FROM toy t
JOIN cat c
    ON t.cat_id = c.cat_id
JOIN owner o
    ON c.owner_id = o.id;
```

toy_name	cat_name	owner_name
ball	Kitty	John Smith
mouse	Kitty	John Smith
ball	Sam	Danielle Davis

JOIN & LEFT JOIN

```
SELECT
    t.toy_name,
    c.cat_name,
    o.name AS owner_name
FROM toy t
JOIN cat c
    ON t.cat_id = c.cat_id
LEFT JOIN owner o
    ON c.owner_id = o.id;
```

toy_name	cat_name	owner_name
ball	Kitty	John Smith
mouse	Kitty	John Smith
ball	Sam	Danielle Davis
mouse	Misty	NULL

LEFT JOIN & LEFT JOIN

```
SELECT
    t.toy_name,
    c.cat_name,
    o.name AS owner_name
FROM toy t
LEFT JOIN cat c
    ON t.cat_id = c.cat_id
LEFT JOIN owner o
    ON c.owner_id = o.id;
```

toy_name	cat_name	owner_name
ball	Kitty	John Smith
mouse	Kitty	John Smith
ball	Sam	Danielle Davis
mouse	Misty	NULL
spring	NULL	NULL

JOIN WITH MULTIPLE CONDITIONS

You can use multiple JOIN conditions using the **ON** keyword once and the **AND** keywords as many times as you need.

CAT AS c				
cat_id	cat_name	mom_id	owner_id	age
1	Kitty	5	1	17
2	Hugo	1	2	10
3	Sam	2	2	5
4	Misty	1	NULL	11

OWNER AS o		
id	name	age
1	John Smith	18
2	Danielle Davis	10

SELECT

```

cat_name,
o.name AS owner_name,
c.age AS cat_age,
o.age AS owner_age
FROM cat c
JOIN owner o
    ON c.owner_id = o.id
    AND c.age < o.age;
```

cat_name	owner_name	age	age
Kitty	John Smith	17	18
Sam	Danielle Davis	5	10

Python For Data Science

Pandas Basics Cheat Sheet

Pandas

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

Use the following import convention:

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

Pandas Data Structures

Series

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

a	3
b	-5
c	7
d	4

Index →

Dataframe

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

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

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

Dropping

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

Asking For Help

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

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
```

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.xlsx')
>>> df = pd.read_excel(xlsx, 'Sheet1')
```

Read and Write to SQL Query or Database Table

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

Selection

Also see NumPy Arrays

Getting

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

Selecting, Boolean Indexing & Setting

By Position

```
>>> df.iloc[[0],[0]] #Select single value by row & column
'Belgium'
>>> df.iat[[0],[0]]
'Belgium'
```

By Label

```
>>> df.loc[[0], ['Country']] #Select single value by row & column labels
'Belgium'
>>> df.at[[0], ['Country']]
'Belgium'
```

By Label/Position

```
>>> df.ix[2] #Select single row of subset of rows
Country Brazil
Capital Brasilia
Population 207847528
>>> df.ix[:,['Capital']] #Select a single column of subset of columns
0 Brussels
1 New Delhi
2 Brasilia
>>> df.ix[1,['Capital']] #Select rows and columns
'New Delhi'
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']>12000000000] #use filter to adjust DataFrame
```

Setting

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

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() #Cumulative 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
>>> 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'])
>>> s3 + 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)
```

~SagaR

Python For Data Science

Matplotlib Cheat Sheet

Matplotlib

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

> Prepare The Data

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'))
```

> 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)
```

> Save Plot

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

> Show Plot

```
>>> plt.show()
```

> 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,colors='blue') #Draw filled polygons
>>> ax.fill_between(x,y,color='yellow') #Fill between y-values and theta
```

2D Data

```
>>> fig, ax = plt.subplots()
>>> im = ax.imshow(img, #Color mapped or RGB arrays
                  cmap='gist_earth',
                  interpolation='nearest',
                  vmin=-2,
                  vmax=2)
>>> 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[1].contourf(data) #Plot filled contours
>>> axes2[2] = ax.clabel(CS) #Label a contour plot
```

Vector Fields

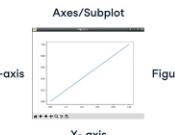
```
>>> axes[0,1].arrow(0,0,0.5,0.5) #Add an arrow to the axes
>>> axes[1,1].quiver(y,z) #Plot a 2D field of arrows
>>> axes[0,1].streamplot(X,Y,U,V) #Plot a 2D field of arrows
```

Data Distributions

```
>>> ax1.hist(y) #Plot a histogram
>>> ax3.boxplot(y) #Make a box and whisker plot
>>> ax3.violinplot(z) #Make a violin plot
```

> Plot Anatomy & Workflow

Plot Anatomy



Workflow

The basic steps to creating plots with matplotlib are:

- 1 Prepare Data
- 2 Create Plot
- 3 Plot
- 4 Customized 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='x')
>>> ax.set_xlim(1, 6.5) #Step 5
>>> plt.savefig('foo.png') #Step 5
>>> plt.show() #Step 6
```

> Close and Clear

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

> Plotting Customize Plot

Colors, Color Bars & Color Maps

```
>>> plt.plot(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(
            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, colors='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 and 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=[3,10,12,'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
```

Python For Data Science

SciPy Cheat Sheet

SciPy



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

> Interacting With NumPy

Also see NumPy

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

Index Tricks

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

Shape Manipulation

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

Polynomials

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

Vectorizing Functions

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

Type Handling

```
>>> np.real(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) #Comb N things taken k time  
>>> misc.central_diff_weights(3) #Weights for N-point central derivative  
>>> misc.derivative(myfunc,1.0) #Find the n-th derivative of a function at a point
```

> Linear Algebra

Also see NumPy

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 #Transpose 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] = 0  
>>> H = sparse.csr_matrix(C) #Compressed Sparse Row matrix  
>>> I = sparse.csr_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
```

Sparse Matrix Decompositions

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

> 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 Functions

```
>>> linalg.sinm(D) #Matrix sine
```

```
>>> linalg.cosm(D) #Matrix cosine
```

```
>>> linalg.tanm(A) #Matrix tangent
```

Hyperbolic Trigonometric Functions

```
>>> linalg.sinhm(D) #Hyperbolic matrix sine
```

```
>>> linalg.coshm(D) #Hyperbolic matrix cosine
```

```
>>> linalg.tanhm(A) #Hyperbolic matrix tangent
```

Matrix Sign Function

```
>>> np.signm(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
```

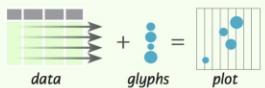
Python For Data Science

Bokeh Cheat Sheet

Plotting With Bokeh

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

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



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

1. Prepare some data (Python lists, NumPy arrays, Pandas DataFrames and other sequences of values)
 2. Create a new plot
 3. Add renderers for your data, with visual customizations
 4. Specify where to generate the output
 5. Show or save the results
- ```
>>> from bokeh.plotting import figure
>>> from bokeh.io import output_file, show
>>> x = [1, 2, 3, 4, 5] #Step 1
>>> y = [6, 7, 2, 4, 5]
>>> p = figure(title="simple line example", #Step 2
X_axis_label='x',
Y_axis_label='y')
>>> p.line(x, y, legend="Temp.", line_width=2) #Step 3
>>> output_file("lines.html") #Step 4
>>> show(p) #Step 5
```

### 1 Data

Also see Lists, NumPy & Pandas

Under the hood, your data is converted to Column Data Sources.

You can also do this manually:

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

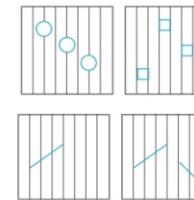
>>> from bokeh.models import ColumnDataSource
>>> cds_df = ColumnDataSource(df)
```

### 2 Plotting

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

## 3 Renderers & Visual Customizations

### Glyphs



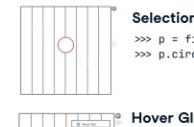
#### Scatter Markers

```
>>> p1.circle(np.array([1,2,3]), np.array([3,2,1]),
fill_color='white')
>>> p2.square(np.array([1,5,3,5,5,5]), [1,4,3],
color='blue', size=1)
```

#### Line Glyphs

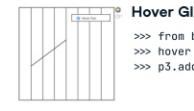
```
>>> p1.line([1,2,3,4], [3,4,5,6], line_width=2)
>>> p2.multi_line(pd.DataFrame([1,2,3],[5,6,7]),
pd.DataFrame([1,3,4,5],[3,2,1])), color='blue')
```

### Customized Glyphs



#### Selection and Non-Selection Glyphs

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



#### Hover Glyphs

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



#### Colormapping

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

### Also see Data

### Legend Location

**Inside Plot Area**

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

#### Outside Plot Area

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

### Legend Orientation

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

### Legend Background & Border

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

### Rows & Columns Layout

#### Rows

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

#### Columns

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

#### Nesting Rows & Columns

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

### Grid Layout

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

### Tabbed Layout

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

### Linked Plots

#### Linked Axes

```
>>> p2.x_range = p1.x_range
>>> p2.y_range = p1.y_range
```

#### Linked Brushing

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

## 5 Show or Save Your Plots

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

# Folium – a *Leaflet* wrapper for

an open-source JavaScript library for mobile-friendly interactive maps



## Quick Start

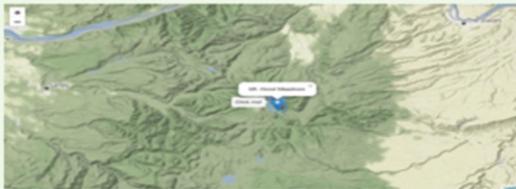
### Installation

Use: `pip install folium`  
or `pip install git+git+goo.gl/kmkGd3` - to install directly from GitHub

### Basic Map

```
> import folium
> from folium import plugins
> m = folium.Map(location = [45.372, -121.6972])
> folium.Marker([45.3288, -121.6625],
 popup='<|>Mt. Hood Meadows</>',
 tooltip='click me').add_to(m)
```

> m



### Markers

#### Icon Markers

> `folium.Marker([lat, lon], popup, tooltip)` - add basic icon markers  
> `folium.Icon(color, icon)` - customize marker icons

#### Circle Markers

> `folium.Circle(location, fill_color, fill_opacity, color, color, _opacity, radius, popup, ...)` - Customize the color, radius, stroke, opacity

#### Vincent/Vega and Altair/VegaLite Markers

> `m = folium.Map(location)`  
> `VegaPop = folium.Popup(max_width).add_child(folium.Vega(json, width, height))` - built-in support for Vincent and altair visualizations  
> `folium.Marker(location, popup=VegaPop).add_to(m)`

#### MarkerClusters

> `m = folium.Map(location)`  
> `marker_cluster = MarkerCluster().add_to(m)`  
> `folium.Marker([lat, lon], popup, tooltip).add_to(marker_cluster)`

#### ClickForMarker and PopOvers

> `m = folium.Map(location)`  
> `m.add_child(folium.LatLngPopup())` - conveniently add lat/long popovers  
> `m.add_child(folium.ClickForMarker(popup))` - on-the-fly placement of markers

#### BoatMarker

> `plugins.BoatMarker(location, heading, wind_heading, wind_speed)`  
- also: color, popup, icon, \*\*kws

## Lines and Shapes

### Lines & PolyLineTextPath

```
> folium.PolyLine(locations, tooltip, popup).add_to(m) - also: color, opacity, weight, smoothing_factor, line_cap, **kws
> plugins.PolyLineTextPath(folium.PolyLine(locations), text, repeat, offset, attributes) - also: orientation, below, center, **kws
```

### Polygons

```
> folium.Polygon(locations, tooltip, popup).add_to(m) - also: color, opacity, weight, fill_color, fill_opacity, smooth_factor, no_clip, **kws
```

### Circles

```
> folium.Circle(location, tooltip, popup).add_to(m) - also: color, opacity, weight, fill_color, fill_opacity, radius, **kws
> folium.CircleMarker(location, tooltip, popup).add_to(m) - also: color, opacity, weight, fill_color, fill_opacity, radius, **kws
```

### Rectangles

```
> folium.Rectangle(bounds=[[lat,lon],[lat,lon]], tooltip, popup).add_to(m) - also: color, fill_color, dash_array, weight, line_join, line_cap, opacity, fill_opacity, **kws
```

## GeoJson and TopoJson

### GeoJson and GeoJsonCSS

```
> folium.GeoJson(GeoJson_path) - add name then use folium.LayerControl().add_to(m)
> folium.GeoJsonCss(GeoJsonCss) - add styling, and popups into the data
```

### TopoJson

```
> folium.TopoJson(TopoJson_path) - add name then use folium.LayerControl().add_to(m)
```

### Choropleth

```
> folium.choropleth(geo_data, data, columns, key_on, fill_color) - also: name, threshold_scale, line_color, line_weight, line_opacity, legend_name, topojson, **kws
```

## Basemaps

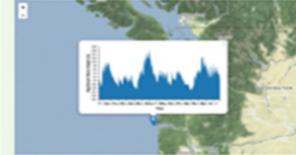
> `folium.Map(location, tiles = 'Stamen Toner')` - use a variety of tiles



## Popups and Labels

> `folium.Popup(html, parse_html, max_width)`

### Vega Popup



### Fancy HTML Popup



## Plugins

### ScrollZoomToggler, Terminator and Fullscreen

```
> plugins.ScrollZoomToggler().add_to(m)
> plugins.Terminator().add_to(m)
> plugins.Fullscreen(position) - also: title, title_cancel, force_serperation_button
```

### Search, Draw and MeasureControl

```
> plugins.Search(GeoJson, search_zoom, geom_type).add_to(m)
- also: search_label, position, popup_on_found
> plugins.Draw().add_to(m) - also: export
> m.add_child(plugins.MeasureControl()) - also: position, primary_length_unit, primary_area_unit, secondary_length_unit, secondary_area_unit
```

### FloatImage

> `plugins.FloatImage(url, bottom, left).add_to(m)`

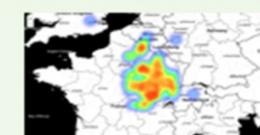
### TimestampedGeoJson

```
> plugins.TimestampedGeoJson(GeoJson, period) - also: add_last_point, dateOptions, auto_play...
```



### HeatMap and HeatMapWithTime

```
> plugins.HeatMap(data).add_to(m) - also: name, radius, min_opacity, max_val, blur, gradient, overlay, max_zoom
> plugins.HeatMapWithTime(data).add_to(m) - also: index, name, radius, min_opacity, max_opacity, auto_play, position, ...
```



# Python For Data Science

## Scikit-Learn Cheat Sheet

### 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=42)
>>> 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', '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)
```

#### > Model Fitting

**Supervised learning**

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

**Unsupervised Learning**

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

#### > Prediction

**Supervised Estimators**

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

**Unsupervised Estimators**

```
>>> y_pred = kmeans.predict(X_test) #Predict labels in clustering algos
```

### > 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
>>> kmeans = KMeans(n_clusters=3, random_state=0)
```

### > Evaluate Your Model's Performance

#### Classification Metrics

**Accuracy Score**

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

**Classification Report**

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

**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<sup>2</sup> 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"]}
>>> search = RandomizedSearchCV(estimator=knn, param_distributions=params,
... cv=4, n_iter=8, random_state=5)
>>> search.fit(X_train, y_train)
>>> print(search.best_score_)
```

# Python For Data Science

## Keras Cheat Sheet

### 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

```
>>> import numpy as np
>>> from keras.models import Sequential
>>> from keras.layers import Dense
>>> data = np.random.randint(2, size=(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)
>>> predictions = model.predict(data)
```

### Data

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

```
>>> 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

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

### Preprocessing

#### Sequence Padding

```
>>> 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

```
>>> 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)
```

### Model Architecture

#### Sequential Model

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

#### Multilayer Perceptron (MLP)

**Binary Classification**

```
>>> 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**

```
>>> 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**

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

#### Convolutional Neural Network (CNN)

```
>>> 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)

```
>>> from keras.layers 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'))
```

### Prediction

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

Also see NumPy & Scikit-Learn

#### Train and Test Sets

```
>>> 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

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

### Inspect Model

```
>>> 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**

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

**MLP: Multi-Class Classification**

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

**MLP: Regression**

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

**Recurrent Neural Network**

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

### Model Training

```
>>> model3.fit(x_train,
 y_train,
 batch_size=32,
 epochs=15,
 verbose=1,
 validation_data=(x_test4,y_test4))
```

### Evaluate Your Model's Performance

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

### Save/ Reload Models

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

### Model Fine-tuning

#### Optimization Parameters

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

#### Early Stopping

```
>>> 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 spaCy Cheat Sheet

## spaCy

spaCy is a free, open-source library for advanced Natural Language processing (NLP) in Python. It's designed specifically for production use and helps you build applications that process and "understand" large volumes of text. Documentation: [spacy.io](https://spacy.io)

```
>>> $ pip install spacy
>>> import spacy
```

## Statistical models

### Download statistical models

Predict part-of-speech tags, dependency labels, named entities and more. See here for available models: [spacy.io/models](https://spacy.io/models)

```
>>> $ python -m spacy download en_core_web_sm
```

### Check that your installed models are up to date

```
>>> $ python -m spacy validate
```

### Loading statistical models

```
>>> import spacy
>>> nlp = spacy.load("en_core_web_sm") # Load the installed model "en_core_web_sm"
```

## Documents and tokens

### Processing text

Processing text with the `nlp` object returns a `Doc` object that holds all information about the tokens, their linguistic features and their relationships

```
>>> doc = nlp("This is a text")
```

### Accessing token attributes

```
>>> doc = nlp("This is a text")
>>> [token.text for token in doc] #Token texts
['This', 'is', 'a', 'text']
```

## Label explanations

```
>>> spacy.explain("RB")
'adverb'
>>> spacy.explain("GPE")
'Countries, cities, states'
```

## > Spans

### Accessing spans

Span indices are exclusive. So `doc[2:4]` is a span starting at token 2, up to – but not including! – token 4.

```
>>> doc = nlp("This is a text")
>>> span = doc[2:4]
>>> span.text
'a text'
```

### Creating a span manually

```
>>> from spacy.tokens import Span #Import the Span object
>>> doc = nlp("I live in New York") #Create a Doc object
>>> span = Span(doc, 3, 5, label="GPE") #Span for "New York" with label GPE (geopolitical)
>>> span.text
'New York'
```

## > Linguistic features

Attributes return label IDs. For string labels, use the attributes with an underscore. For example, `token.pos_`.

### Part-of-speech tags

Predicted by Statistical model

```
>>> doc = nlp("This is a text.")
>>> [token.pos_ for token in doc] #Coarse-grained part-of-speech tags
['DET', 'VERB', 'DET', 'NOUN', 'PUNCT']
>>> [token.tag_ for token in doc] #Fine-grained part-of-speech tags
['DT', 'VBD', 'DT', 'NN', '.']
```

### Syntactic dependencies

Predicted by Statistical model

```
>>> doc = nlp("This is a text.")
>>> [token.dep_ for token in doc] #Dependency labels
['nsubj', 'ROOT', 'det', 'attr', 'punct']
>>> [token.head.text for token in doc] #Syntactic head token (governor)
['is', 'is', 'text', 'is', 'is']
```

### Named entities

Predicted by Statistical model

```
>>> doc = nlp("Larry Page founded Google")
>>> [(ent.text, ent.label_) for ent in doc.ents] #Text and label of named entity span
[('Larry Page', 'PERSON'), ('Google', 'ORG')]
```

## > Pipeline components

Functions that take a `Doc` object, modify it and return it.



### Pipeline information

```
>>> nlp = spacy.load("en_core_web_sm")
>>> nlp.pipe_names
['tagger', 'parser', 'ner']
>>> nlp.pipeline
[(tagger, <spacy.pipeline.Tagger>),
(parser, <spacy.pipeline.DependencyParser>),
(ner, <spacy.pipeline.EntityRecognizer>)]
```

### Custom components

```
def custom_component(doc): #Function that modifies the doc and returns it
 print("Do something to the doc here!")
 return doc
nlp.add_pipe(custom_component, first=True) #Add the component first in the pipeline
Components can be added first, last (default), or before or after an existing component.
```

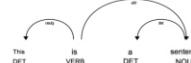
## > Visualizing

If you're in a Jupyter notebook, use `displacy.render` otherwise, use `displacy.serve` to start a web server and show the visualization in your browser.

```
>>> from spacy import displacy
```

### Visualize dependencies

```
>>> doc = nlp("This is a sentence")
>>> displacy.render(doc, style="dep")
```



### Visualize named entities

```
>>> doc = nlp("Larry Page founded Google")
>>> displacy.render(doc, style="ent")
```

```
Larry Page PERSON founded Google ORG
```

## > Word vectors and similarity

To use word vectors, you need to install the larger models ending in `md` or `lg`, for example `en_core_web_lg`.

### Comparing similarity

```
>>> doc1 = nlp("I like cats")
>>> doc2 = nlp("I like dogs")
>>> doc1.similarity(doc2) #Compare 2 documents
>>> doc1[2].similarity(doc2[2]) #Compare 2 tokens
>>> doc1[0].similarity(doc2[1:3]) # Compare tokens and spans
```

### Accessing word vectors

```
>>> doc = nlp("I like cats") #Vector as a numpy array
>>> doc[2].vector #The L2 norm of the token's vector
>>> doc[2].vector_norm
```

## > Syntax iterators

### Sentences

Usually needs the dependency parser

```
>>> doc = nlp("This a sentence. This is another one.")
>>> [sent.text for sent in doc.sents] #doc.sents is a generator that yields sentence spans
['This is a sentence.', 'This is another one.')
```

### Base noun phrases

Needs the tagger and parser

```
>>> doc = nlp("I have a red car")
#doc.noun_chunks is a generator that yields spans
>>> [chunk.text for chunk in doc.noun_chunks]
['I', 'a red car']
```

~SagaR

# Python For Data Science

## PySpark SQL Basics Cheat Sheet

### PySpark & Spark SQL



Spark SQL is Apache Spark's module for working with structured data.

#### > Initializing SparkSession

A SparkSession can be used create DataFrame, register DataFrame as tables, execute SQL over tables, cache tables, and read parquet files.

```
>>> from pyspark.sql import SparkSession
>>> spark = SparkSession \
 .builder \
 .appName("Python Spark SQL basic example") \
 .config("spark.some.config.option", "some-value") \
 .getOrCreate()
```

#### > Creating DataFrames

##### From RDDs

```
>>> from pyspark.sql.types import *
Infer Schema
>>> sc = spark.sparkContext
>>> lines = sc.textFile("people.txt")
>>> parts = lines.map(lambda l: l.split(","))
>>> people = parts.map(lambda p: Row(name=p[0],age=int(p[1])))
>>> peopleDF = spark.createDataFrame(people)
Specify Schema
>>> people = parts.map(lambda p: Row(name=p[0],
 age=int(p[1].strip())))
+-----+-----+
| name | age |
+-----+-----+
| Mine | 28 |
| Filip | 29 |
| Jonathan | 30 |
+-----+-----+
```

##### From Spark Data Sources

**JSON**

```
>>> df = spark.read.json("customer.json")
>>> df.show()
+-----+-----+
| address|lastName|name|phoneNumber|
+-----+-----+
| New York,10021,N...|John|Smith|1234 555-1234,ho...
| New York,10021,N...|John|Smith|1234 555-1234,ho...
| New York,10021,N...|John|Smith|1234 555-1234,ho...
+-----+-----+
```

**Parquet files**

```
>>> df3 = spark.read.load("people.json", format="json")
```

**TXT files**

```
>>> df4 = spark.read.text("people.txt")
```

#### > Filter

```
#Filter entries of age, only keep those records of which the values are >24
>>> df.filter(df["age"]>24).show()
```

### > Duplicate Values

```
>>> df = df.dropDuplicates()
```

### > Queries

```
>>> from pyspark.sql import functions as F
Select
>>> df.select("firstName").show() #Show all entries in firstName column
>>> df.select("firstName", "lastName") \
 .show()
>>> df.select("firstName", #Show all entries in firstName, age and type
 "age",
 explode("phoneNumber") \
 .alias("contactInfo")) \
 .select("contactInfo.type",
 "firstName",
 "lastName",
 "age") \
 .show()
>>> df.select(df["firstName"],df["age"]+1) #Show all entries in firstName and age,
 .show() #add 1 to the entries of age
>>> df.select(df['age'] > 24).show() #Show all entries where age >24
```

```
When
>>> df.select("firstName", #Show firstName and 0 or 1 depending on age >30
 F.when(df.age > 30, 1) \
 .otherwise(0)) \
 .show()
```

```
>>> df[df.firstName.isin("Jane","Boris")] #Show firstName if in the given options
 .collect()
```

#### Like

```
>>> df.select("firstName", #Show firstName, and lastName is TRUE if lastName is like Smith
 df.lastName.like("Smith")) \
 .show()
```

#### Startswith - Endswith

```
>>> df.select("firstName", #Show firstName, and TRUE if lastName starts with Sm
 df.lastName \
 .startswith("Sm")) \
 .show()
```

```
>>> df.select(df.lastName.endswith("th")) #Show last names ending in th
 .show()
```

#### Substring

```
>>> df.select(df.firstName.substr(1, 3) \ #Return substrings of firstName
 .alias("name")) \
 .collect()
```

#### Between

```
>>> df.select(df.age.between(22, 24)) \ #Show age: values are TRUE if between 22 and 24
 .show()
```

### > Add, Update & Remove Columns

#### Adding Columns

```
>>> df = df.withColumn('city',df.address.city) \
 .withColumn('postalCode',df.address.postalCode) \
 .withColumn('state',df.address.state) \
 .withColumn('streetAddress',df.address.streetAddress) \
 .withColumn('telephoneNumber', explode(df.phoneNumber.number)) \
 .withColumn('phoneType', explode(df.phoneNumber.type))
```

#### Updating Columns

```
>>> df = df.withColumnRenamed('telephoneNumber', 'phoneNumber')
```

#### Removing Columns

```
>>> df = df.drop("address", "phoneNumber")
>>> df = df.drop(df.address).drop(df.phoneNumber)
```

### > Missing & Replacing Values

```
>>> df.na.fill(50).show() #Replace null values
>>> df.na.drop().show() #Return new df omitting rows with null values
>>> df.na \ #Return new df replacing one value with another
 .replace(10, 20) \
 .show()
```

### > GroupBy

```
>>> df.groupBy("age")\ #Group by age, count the members in the groups
 .count() \
 .show()
```

### > Sort

```
>>> peopleDF.sort(peopleDF.age.desc()).collect()
>>> df.sort("age", ascending=False).collect()
>>> df.orderBy([("age", "city"], ascending=[0,1])\
 .collect()
```

### > Repartitioning

```
>>> df.repartition(10)\ #df with 10 partitions
 .rdd \
 .getNumPartitions()
>>> df.coalesce(1).rdd.getNumPartitions() #df with 1 partition
```

### > Running Queries Programmatically

#### Registering DataFrames as Views

```
>>> peopleDF.createGlobalTempView("people")
>>> df.createTempView("customer")
>>> df.createOrReplaceTempView("customer")
```

#### Query Views

```
>>> df5 = spark.sql("SELECT * FROM customer").show()
>>> peopleDF2 = spark.sql("SELECT * FROM global_temp.people")\
 .show()
```

### > Inspect Data

```
>>> df.dtypes #Return df column names and data types
>>> df.show() #Display the content of df
>>> df.head() #Return first n rows
>>> df.first() #Return first row
>>> df.take(2) #Return the first n rows >> df.schema Return the schema of df
>>> df.describe().show() #Compute summary statistics >> df.columns Return the columns of df
>>> df.count() #Count the number of rows in df
>>> df.distinct().count() #Count the number of distinct rows in df
>>> df.printSchema() #Print the schema of df
>>> df.explain() #Print the (logical and physical) plans
```

### > Output

#### Data Structures

```
>>> rdd = df.rdd #Convert df into an RDD
>>> df.toJSON().first() #Convert df into a RDD of string
>>> df.toPandas() #Return the contents of df as Pandas DataFrame
```

#### Write & Save to Files

```
>>> df.select("firstName", "city")\
 .write \
 .save("nameAndCity.parquet")
>>> df.select("firstName", "age") \
 .write \
 .save("namesAndAges.json",format="json")
```

### > Stopping SparkSession

```
>>> spark.stop()
```

~SagaR

# Python For Data Science

## PySpark RDD Cheat Sheet

### 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
```

Select 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 any Hadoop-supported file system URI with `textFile()`, or read in a directory of text files with `wholeTextFiles()`

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

### > Retrieving RDD Information

#### Basic Information

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

#### Summary

```
>>> rdd3.max() #Maximum value of RDD elements
99
>>> rdd3.min() #Minimum value of RDD elements
0
>>> rdd3.mean() #Mean value of RDD elements
49.5
>>> rdd3.stdev() #Standard deviation of RDD elements
28.86607904772218
>>> rdd3.variance() #Compute variance of RDD elements
833.25
>>> rdd3.histogram(3) #Compute histogram by bins
([0,33,66,99],[33,33,34])
>>> rdd3.stats() #Summary statistics (count, mean, stdev, max & min)
```

### > Applying Functions

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

### > Selecting Data

**Getting**

```
>>> rdd.collect() #Return a list with all RDD elements
[('a', 7), ('a', 2), ('b', 2)]
>>> rdd.take(2) #Take first 2 RDD elements
[('a', 7), ('a', 2)]
>>> rdd.first() #Take first RDD element
('a', 7)
>>> rdd.top(2) #Take top 2 RDD elements
[('b', 2), ('a', 7)]
Sampling
>>> rdd3.sample(False, 0.15, 81).collect() #Return sampled subset of rdd3
[3,4,27,31,48,41,42,43,60,76,79,80,86,97]
Filtering
>>> rdd.filter(lambda x: "a" in x).collect() #Filter the RDD
[('a',7),('a',2)]
>>> rdd5.distinct().collect() #Return distinct RDD values
[('a',2),('b',7)]
>>> rdd.keys().collect() #Return (key,value) RDD's keys
['a', 'a', 'b']
```

### > Iterating

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

### > Reshaping Data

#### Reducing

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

#### Grouping by

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

#### Aggregating

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

### > Mathematical Operations

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

### > Sort

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

### > 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
```

~SagaR

| Import Resources                                         | Make a soup object out of a website                                                                                                                                                                                                                                                                                                 | Object types                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                |
|----------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <pre>import requests from bs4 import BeautifulSoup</pre> | <pre>// 1. The HTTP request webpage = request.get('URL', 'html.parser');  // 2. Turn the website into a soup object soup = BeautifulSoup(webpage.content);</pre> <p>"html.parser" is one option for parsers we could use. There are other options, like "lxml" and "html5lib" that have different advantages and disadvantages.</p> | <pre>//1. Tags correspond to HTML tags <b>Example Code:</b> soup = BeautifulSoup('&lt;div id="example"&gt;An example div&lt;/div&gt;&lt;p&gt;An example p tag&lt;/p&gt;');  print(soup.div); --&gt; &lt;div id="example"&gt;An example div&lt;/div&gt; --&gt; gets the <b>first</b> tag of that type on the page  print(soup.div.name) print(soup.div.attrs) --&gt; div --&gt; {'id': 'example'}</pre><br><pre>//2. Navigable Strings: Piece of text inside of HTML Tags print(soup.div.string) --&gt; An example div</pre> |

| Navigating by Tasks                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | Finding All                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          | CSS Selectors                                                                                                                                                                                                                                                                                         |
|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p><b>Example Code:</b></p> <pre>&lt;h1&gt;World's Best Chocolate Chip Cookies&lt;/h1&gt; &lt;div class="banner"&gt; &lt;h1&gt;Ingredients&lt;/h1&gt; &lt;/div&gt; &lt;ul&gt; &lt;li&gt; 1 cup flour &lt;/li&gt; &lt;li&gt; 1/2 cup sugar &lt;/li&gt; &lt;li&gt; 2 tbsp oil &lt;/li&gt; &lt;li&gt; 1/2 tsp baking soda &lt;/li&gt; &lt;li&gt; 1/2 cup chocolate chips &lt;/li&gt; &lt;li&gt; 1/2 tsp vanilla &lt;/li&gt; &lt;li&gt; 2 tbsp milk &lt;/li&gt; &lt;/ul&gt;  //1. Get the children of a tag: for child in soup.ul.children:     print(child) --&gt; &lt;li&gt; 1 cup flour &lt;/li&gt; --&gt; &lt;li&gt; 1/2 cup sugar &lt;/li&gt; ... //2. Get the parent of a tag: for parent in soup.li.parents:     print(parent)</pre> | <pre>//1. <b>find_all()</b> print(soup.find_all("h1")) --&gt; Outputs all &lt;h1&gt;...&lt;/h1&gt; on the website  //1.1. <b>find_all()</b> with <b>regex</b> import re soup.find_all(re.compile("[ou]l")) --&gt; Outputs all &lt;ul&gt;...&lt;/ul&gt; and &lt;ol&gt;...&lt;/ol&gt; soup.find_all(re.compile("h[1-9]")) --&gt; Outputs all headings  //1.2. <b>find_all()</b> with <b>lists</b> soup.find_all(['h1', 'a', 'p'])  //1.3 <b>find_all()</b> with <b>attributes</b> soup.find_all(attrs={'class':'banner', 'id':'jumbotron'});  //1.4 <b>find_all()</b> with <b>functions</b> def has_banner_class_and_hello_world(tag):     return tag.attr('class') == "banner" and tag.string == "Hello world"  soup.find_all(has_banner_class_and_hello_world)</pre> | <pre>//1. grab CSS classes with .select("class_name") soup.select(".recipeLink")  //*2. grab CSS IDs with .select("#id_name") soup.select("#selected")  //3. using a loop for link in soup.select(".recipeLink &gt; a"):     webpage = requests.get(link)     new_soup = BeautifulSoup(webpage)</pre> |

**BeautifulSoup**  
**~SagaR**

# Selenium Python Cheat Sheet

| 1. Import the Selenium library                                                                                                                                                                                                                                                           | 2. Driver Initialization with Python                                                                                                                                                                                                                                          | 3. Path to uses the download file                                                                                                                                                     |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| from selenium import webdriver                                                                                                                                                                                                                                                           | # Initialize Chrome WebDriver<br>driver = webdriver.Chrome()                                                                                                                                                                                                                  | driver =<br>webdriver.Chrome(executable_path="/usr/bin/chromedriver/",<br>options=chrome_options )                                                                                    |
| 3. Setting Options in Selenium WebDriver                                                                                                                                                                                                                                                 | 4. Finding an element                                                                                                                                                                                                                                                         | 5. Locate Element by CSS Class                                                                                                                                                        |
| from selenium import webdriver<br><br>from selenium.webdriver.chrome.options import Options                                                                                                                                                                                              | from selenium import webdriver<br><br>driver = webdriver.Chrome()<br>driver.get("https://www.lambdatest.com")<br>email_form =<br>driver.find_element_by_id("testing_form")                                                                                                    | from selenium import webdriver<br><br>driver = webdriver.Chrome()<br>driver.get("https://www.lambdatest.com")<br>first_form_input = driver.find_element_by_class_name("form-control") |
| 6. Locate Elements by Name                                                                                                                                                                                                                                                               | 7. Locate Elements by XPath                                                                                                                                                                                                                                                   | 8. Locating element by tag name                                                                                                                                                       |
| Ex: <input name="name" type="text" value="user name" /><br><br>from selenium import webdriver<br><br>driver = webdriver.Chrome()<br>driver.get("https://www.lambdatest.com")<br><br># for selection input with name attribute "name"<br>name_input = driver.find_element_by_name("name") | from selenium import webdriver<br><br>driver = webdriver.Chrome()<br><br>driver.get("https://www.lambdatest.com")<br><br>email_input =<br>driver.find_element_by_xpath("//input[@name='email']")                                                                              | driver = webdriver.Chrome()<br><br>driver.get("https://www.lambdatest.com")<br><br>email_input =<br>driver.find_element_by_tag_name("input")                                          |
| 9. Locate Element by Link text or Partial Link Text                                                                                                                                                                                                                                      | 10. Misc methods for finding elements                                                                                                                                                                                                                                         | 11. Opening a URL (or document)                                                                                                                                                       |
| from selenium import webdriver<br><br>driver = webdriver.Chrome()<br>driver.get("https://www.lambdatest.com")<br><br>email_input = driver.find_element_by_link_text('Start Free Testing')                                                                                                | <ul style="list-style-type: none"> <li>□ ID = "id"</li> <li>□ XPATH = "xpath"</li> <li>□ NAME = "name"</li> <li>□ TAG_NAME = "tag name"</li> <li>□ CLASS_NAME = "class name"</li> <li>□ LINK_TEXT = "link text"</li> <li>□ PARTIAL_LINK_TEXT = "partial link text"</li> </ul> | driver.get('url')<br><br>Refresh Page<br>driver.refresh()                                                                                                                             |

|                                                                                                                                                                                                                  |                                                                                                                                              |                                                                                                                                                                                                                                                    |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>12. Clear in text box</b>                                                                                                                                                                                     | <b>13. Click on object</b>                                                                                                                   | <b>14. Dragging and Dropping a WebElement</b>                                                                                                                                                                                                      |
| <pre># get element element = driver.find_element_by_id("useremail")  # send keys element.clear()</pre>                                                                                                           | <pre># get element button_element = driver.find_element_by_link_text("Start Free Testing")  # click the element button_element.click()</pre> | <pre>element = driver.find_element_by_name("source") target = driver.find_element_by_name("target")  from selenium.webdriver import ActionChains action_chains = ActionChains(driver) action_chains.drag_and_drop(element, target).perform()</pre> |
| <b>15. Selecting an option</b>                                                                                                                                                                                   | <b>16. Navigating between windows</b>                                                                                                        | <b>17. driver.window_handles</b>                                                                                                                                                                                                                   |
| <pre>from selenium.webdriver.support.ui import Select select = Select(driver.find_element_by_id('city')) select.select_by_index(index) select.select_by_visible_text("text") select.select_by_value(value)</pre> | driver.switch_to_window("window_handle")                                                                                                     | for handle in driver.window_handles:     driver.switch_to_window(handle)                                                                                                                                                                           |
| <b>18. driver.current_window_handle</b>                                                                                                                                                                          | <b>19. Switching to iFrames</b>                                                                                                              | <b>20. driver.switch_to_default_content()</b>                                                                                                                                                                                                      |
| handler = driver.current_window_handle                                                                                                                                                                           | driver.switch_to_frame("frame_name.0.child")                                                                                                 | driver.switch_to_default_content()                                                                                                                                                                                                                 |
| <b>21. Handling pop-ups and alerts</b>                                                                                                                                                                           | <p>→ <u>alert_obj.text()</u></p> <pre>alert_obj = driver.switch_to.alert msg = alert_obj.text() print(msg)</pre>                             | <b>22. Getting Page Source</b>                                                                                                                                                                                                                     |
| <input type="checkbox"/> Simple Alert<br><input type="checkbox"/> Confirmation Alert<br><input type="checkbox"/> Prompt Alert                                                                                    |                                                                                                                                              | page_source = driver.page_source                                                                                                                                                                                                                   |
| <b>→ <u>driver.switch_to.alert</u></b><br><pre>alert_obj = driver.switch_to.alert</pre>                                                                                                                          |                                                                                                                                              | <b>23. Navigating browser history</b>                                                                                                                                                                                                              |
| <b>→ <u>alert_obj.accept()</u></b><br><pre>alert_obj = driver.switch_to.alert alert_obj.accept()</pre>                                                                                                           |                                                                                                                                              | driver.forward()                                                                                                                                                                                                                                   |
| <b>→ <u>alert_obj.dismiss()</u></b><br><pre>alert_obj = driver.switch_to.alert alert_obj.accept()</pre>                                                                                                          |                                                                                                                                              |                                                                                                                                                                                                                                                    |

# Python For Data Science

## Data Wrangling in Pandas Cheat Sheet

### > Reshaping Data

#### Pivot

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

#### Pivot Table

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

#### Stack / Unstack

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

#### Melt

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

### > Iteration

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

### > Missing Data

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

### > Advanced Indexing

Also see NumPy Arrays

#### Selecting

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

#### Indexing With isin()

```
>>> df[(df.Country.isin(df2.Type))] #Find same elements
>>> df3.filter(items="a","b")) #Filter on values
>>> df.select(lambda x: not x%5) #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') #Set the index
>>> df4 = df.reset_index() #Reset the index
>>> df = df.rename(index=str, #Rename
 DataFrame columns={'Country':'cntry',
 'Capital':'cptl',
 'Population':'pltn'})
```

### Reindexing

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

#### Forward Filling

```
>>> df.reindex(range(4),
 method='ffill')
```

| Country   | Capital   | Population | 0 | 3 |
|-----------|-----------|------------|---|---|
| 0 Belgium | Brussels  | 11190846   | 1 | 3 |
| 1 India   | New Delhi | 1303171035 | 2 | 3 |
| 2 Brazil  | Brasilia  | 207847528  | 3 | 3 |
| 3 Brazil  | Brasilia  | 207847528  | 4 | 3 |

#### Backward Filling

```
>>> s3 = s.reindex(range(5),
 method='bfill')
```

| 0 | 3 |
|---|---|
| 1 | 3 |
| 2 | 3 |
| 3 | 3 |
| 4 | 3 |

### Multilabeling

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

### > Combining Data

| data1 | data2  |
|-------|--------|
| X1    | X2     |
| a     | 11.432 |
| b     | 1.303  |
| c     | 99.906 |

#### 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 |

#### 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')
```

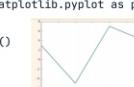
### > Dates

```
>>> 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

```
>>> import matplotlib.pyplot as plt
>>> s.plot()
>>> plt.show()
```

Also see Matplotlib



~SagAR

# Python For Data Science

## Importing Data Cheat Sheet

### > 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

```
>>> 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)
```

##### 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 type

```
>>> 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
 comments='#', #Character to split comments
 na_values=['']) #String to recognize as NA/NaN
```

### > 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 NumPy array
```

### > SAS File

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

### > Stata File

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

### > Excel Spreadsheets

```
>>> file = 'urbanpop.xlsx'
>>> data = pd.ExcelFile(file)
>>> df_sheet2 = data.parse('1960-1964',
 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
```

### > 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)
```

### > Pickled Files

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

### > Matlab Files

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

### > HDF5 Files

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

### > Exploring Dictionaries

Querying relational databases with pandas

```
>>> print(mat.keys()) #Print dictionary keys
>>> for key in mat.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
#Retrieve the value for a key
>>> print(data['meta']['Description'].value)
```

### > 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", "test2.txt") #Rename a file
>>> os.remove("test1.txt") #Delete an existing file
>>> os.mkdir("newdir") #Create a new directory
```

```
|!pip install pandas
```

### 13. Read SQL Table

`import pandas as pd`  
We can extract table from SQL database (SQL Server / Teradata). See the program below -

We can load Stata data file via `read_stata()` function.  
`mydata41 = pd.read_stata('cars.dta')`

#### SQL Server

You can read data from tables stored in SQL Server by building a connection. You need to have server, User ID (UID), database details to establish connection.

```
import pandas as pd
import pyodbc
conn = pyodbc.connect("Driver={SQL Server};Server=serverName;UID=UserName;PWD=Password;Database=RCO_DW;")
df = pd.read_sql_query('select * from dbo.Table WHERE ID > 10', conn)
df.head()
```

#### Teradata

You need to import **Teradata module** which makes python easily integrated with Teradata Database. [Install this package](#) using the command below -

```
import pandas as pd
import teradata
udaExec = teradata.UdaExec(appName="HelloWorld", version="1.0",
 logConsole=False)
session = udaExec.connect(method="odbc",
 USERREGIONALSETTINGS="N",
 system="tdprod",
 username="xxx",
 password="xxx");
```

```
query = "SELECT * FROM flight"
df = pd.read_sql(query, session)
```

#### Explanation

**4. Read Excel** ■ UdaExec provides DevOps support features such as configuration and logging.

`mydata = pd.read_excel('flight.xlsx', skiprows=2, engine='odbc', header=0, sheet_name='Sheet1')` ■ You can assign any name and version in `appName` and `version`

`skiprows=2` ■ `logConsole=False` tells Python not to log to the console.

`If you do not specify system` ■ `system="tdprod"` refers to name of the system we are connecting using ODBC as the connection method

`USERREGIONALSETTINGS="N"` ■ `USERREGIONALSETTINGS="N"` is used to ensure that float values can be loaded and make decimal separator be ','

Suppose you have .db extension file which is a database file and you want to extract data from it. | [11. Specify values as missing values](#)

```
import sqlite3
from pandas.io import sql
conn = sqlite3.connect('C:/Users/Deepanshu/Downloads/flight.db')
query = "SELECT * FROM flight"
results = pd.read_sql(query, con=conn)
print results.head()
```

### 6. Read SAS File

We can import SAS data file by using `read_sas()` function.

```
mydata4 = pd.read_sas('cars.sas7bdat')
```

### 12. Skip rows while importing

Suppose you want to skip first 5 rows and wants to read data from 6th row (6th row would be a header row)

```
mydata8 = pd.read_csv("http://winterolympicsmedals.com/medals.csv", skiprows=5)
```

