

Machine Learning & Intelligence Artificielle



 cpc
Manuel Simoes
manuel.simoes@cpc-analytics.fr

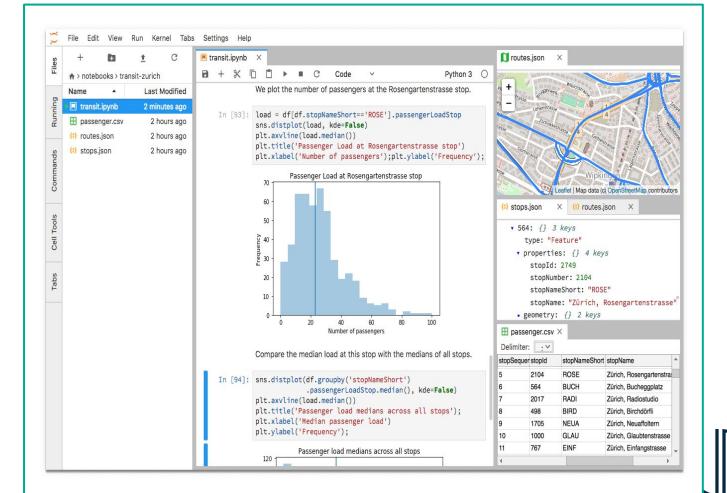
Les outils



- Utilisation du langage Python et de ses librairies
 - Si possible, installez également CUDA (surtout pour les réseaux de neurones)

- Installez
 - Anaconda 2020.07 avec Python 3.8
 - Gère de nombreuses librairies : Pandas, Scikit-learn, Numpy, Scipy, Matplotlib, ...
 - Keras et Tensor Flow (réseaux de neurones compatibles avec GPU)

- Travail sur Notebook
 - Utilisez Jupyter lab / Jupyter notebook

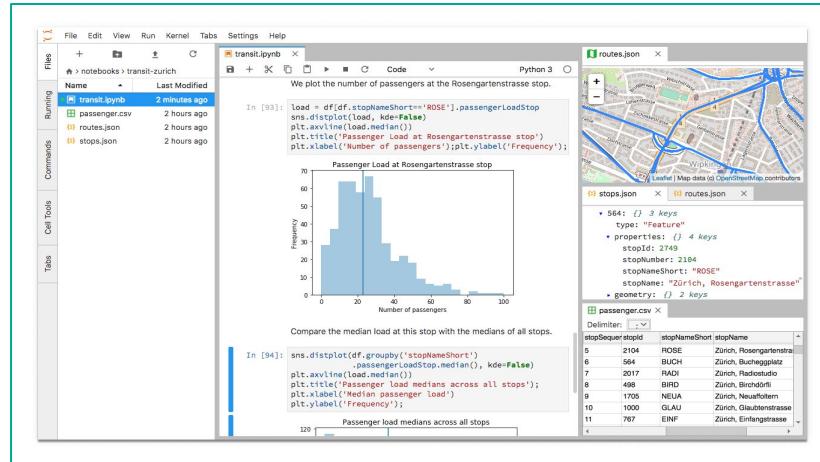
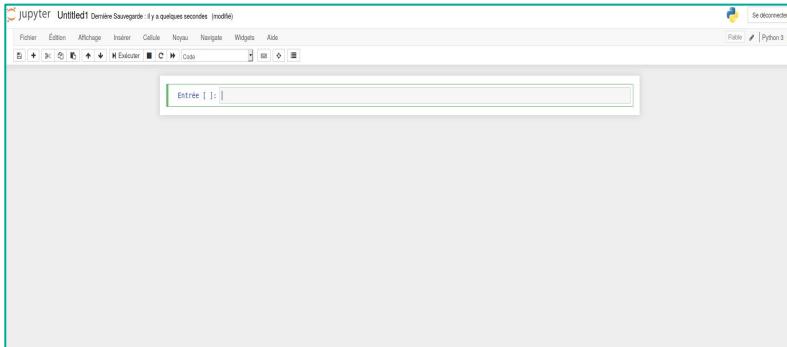


- Premier outils -

Les modules Python utile pour le Machine Learning



NoteBook



Changer les tailles maximales des tableaux (ainsi que le nombre de ligne maximum) affichés dans les output du Notebook.

```
pd.set_option('display.max_columns', 200)  
pd.set_option('display.max_rows', 500)
```

Pour utiliser toute la largeur de la fenêtre du navigateur.

```
from IPython.display import HTML  
display(HTML("<style>.container { width: 100% !important; }</style>"))
```

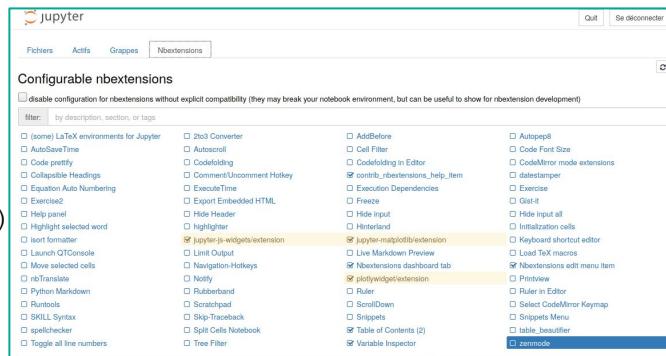
Touche rapide pratique :

Pour créer une cellule avant la cellule que vous avez sélectionnée (elle est encadrée en bleu pas en vert) : *a*
ou après la cellule sélectionnée : *b*

Supprimer la cellule sélectionnée : *dd*

Imprimer dans la cellule des informations qui vous intéresse :

```
display("text ou instruction")
```



https://github.com/ipython-contrib/jupyter_contrib_nbextensions

Manipulation de données

Les librairies les plus importantes et incontournables pour le travail avec des données. D'autres librairies permet d'utiliser leur syntaxe tout en utilisant les ressources GPU le plus souvent possibles.



The fundamental package for scientific computing with Python

GET STARTED

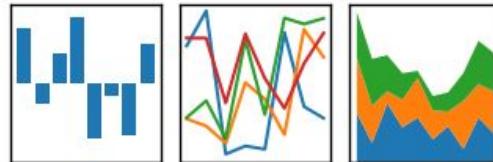
<https://numpy.org/>



<https://www.scipy.org/>

pandas

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



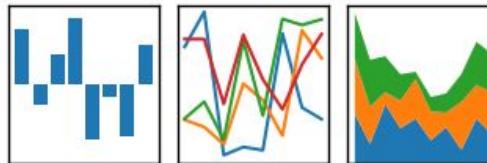
<https://pandas.pydata.org/>



Modules utile pour la Data Science pour Python

pandas

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



<https://pandas.pydata.org/>

scikit-learn

Machine Learning in Python

[Getting Started](#) [Release Highlights for 0.23](#) [GitHub](#)

- Simple and efficient tools for predictive data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

<https://scikit-learn.org/stable/index.html>

NumPy



The fundamental package for scientific computing with Python

[GET STARTED](#)

<https://numpy.org/>

SciPy.org



<https://www.scipy.org/>

Réseau de neurone



Keras

Simple. Flexible. Powerful.

<https://keras.io/>



TensorFlow

<https://www.tensorflow.org/>



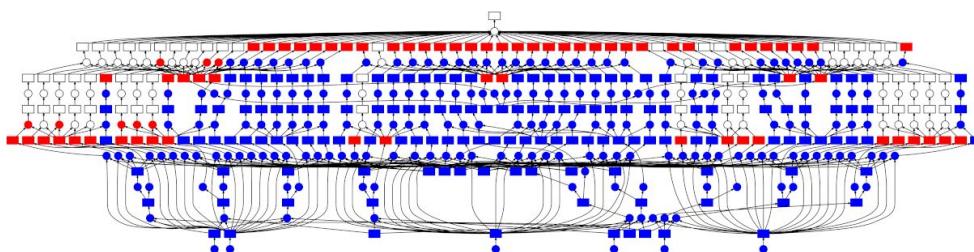
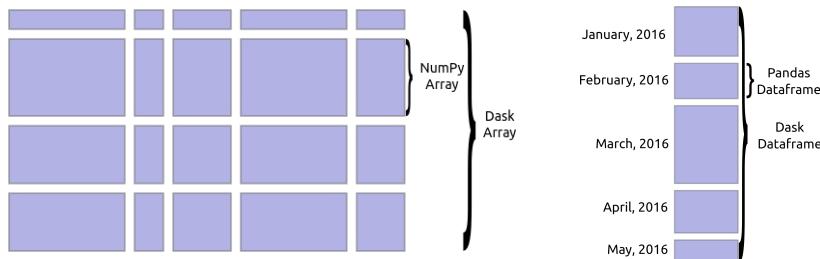
PyTorch

<https://pytorch.org/>

Puissance de calcul



<https://dask.org/>



Autre Modules python pour la Data Science



statsmodels is a Python module that provides classes and functions for the estimation of many different statistical models, as well as for conducting statistical tests, and statistical data exploration. An extensive list of result statistics are available for each estimator. The results are tested against existing statistical packages to ensure that they are correct. The package is released under the open source Modified BSD (3-clause) license. The online documentation is hosted at [statsmodels.org](https://www.statsmodels.org).

<https://www.statsmodels.org/stable/index.html>

Python Outlier Detection (PyOD)

Deployment & Documentation & Stats

[pypi v0.8.4](#) [Anaconda Cloud 0.8.4](#) [docs passing](#) [launch binder](#) [stars 3.8k](#) [forks 787](#) [downloads 1M](#) [downloads/month 93k](#)

Build Status & Coverage & Maintainability & License

[build passing](#) [build passing](#) [PASSED](#) [coverage 95%](#) [maintainability](#) [license BSD-2-Clause](#)

PyOD is a comprehensive and scalable **Python toolkit** for **detecting outlying objects** in multivariate data. This exciting yet challenging field is commonly referred as [Outlier Detection](#) or [Anomaly Detection](#).

PyOD includes more than 30 detection algorithms, from classical LOF (SIGMOD 2000) to the latest COPOD (ICDM 2020). Since 2017, PyOD has been successfully used in numerous academic researches and commercial products [9] [18] [29] [31]. It is also well acknowledged by the machine learning community with various dedicated posts/tutorials, including [Analytics Vidhya](#), [KDnuggets](#), [Towards Data Science](#), [Computer Vision News](#), and [awesome-machine-learning](#).

PyOD is featured for:

- **Unified APIs, detailed documentation, and interactive examples** across various algorithms.
- **Advanced models**, including **classical ones from scikit-learn, latest deep learning methods**, and emerging algorithms like COPOD.
- **Optimized performance with JIT and parallelization** when possible, using [numba](#) and [joblib](#).
- **Compatible with both Python 2 & 3**.

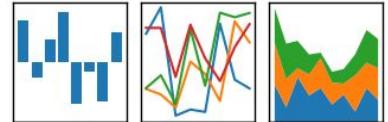
<https://github.com/yzhao062/Pyod>

Voir aussi https://scikit-learn.org/stable/modules/outlier_detection.html

Visualisation 1 / 10

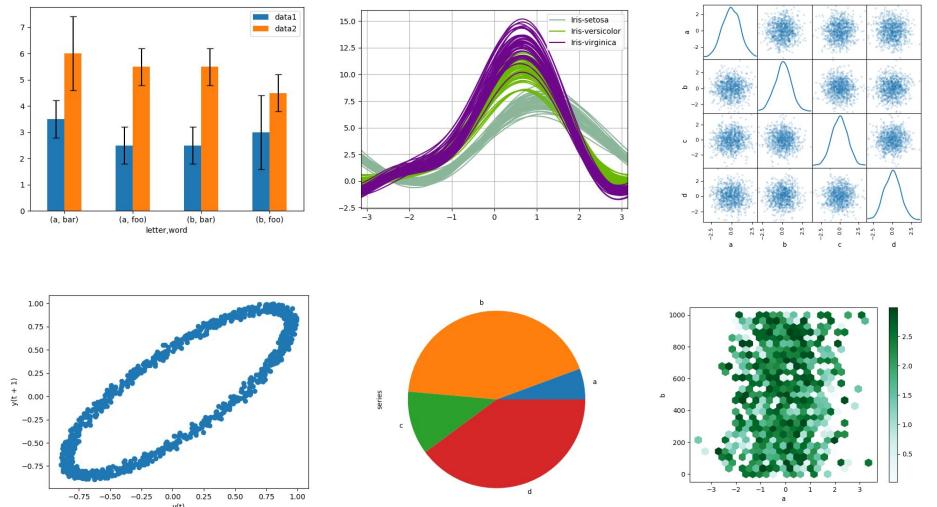
pandas

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



Pandas - visualization:

Pandas contient des méthodes qui permettent de représenter très rapidement un graphe (histogramme,...).



Basic plotting: `plot`

We will demonstrate the basics, see the [cookbook](#) for some advanced strategies.

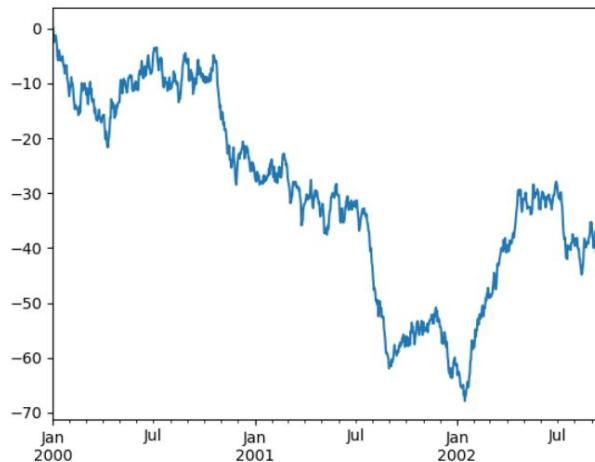
The `plot` method on Series and DataFrame is just a simple wrapper around `plt.plot()`:

```
In [3]: ts = pd.Series(np.random.randn(1000),
...:                     index=pd.date_range('1/1/2000', periods=1000))
...:

In [4]: ts = ts.cumsum()

In [5]: ts.plot()
Out[5]: <matplotlib.axes._subplots.AxesSubplot at 0x7f4505e618d0>
```

https://pandas.pydata.org/pandas-docs/stable/user_guide/visualization.html



Visualisation 2 / 10



<https://matplotlib.org/>

matplotlib

C'est la librairie graphique de référence de Python pour créer des figures et des graphes.

Elle est utilisé par Pandas et elle permet d'adapter les figures de Pandas. Elle n'est pas intuitive et demande souvent plusieurs lignes de code.

Gallery

This gallery contains examples of the many things you can do with Matplotlib. Click on any image to see the full image and source code.

For longer tutorials, see our [tutorials page](#). You can also find [external resources](#) and a [FAQ](#) in our [user guide](#).

Lines, bars and markers

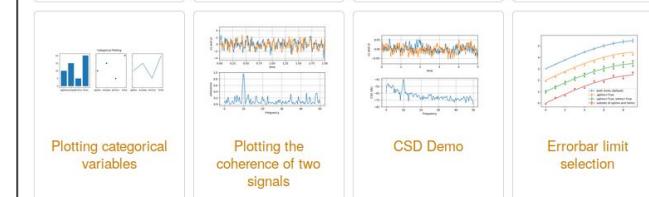


Stacked Bar Graph

Grouped bar chart with labels

Horizontal bar chart

Broken Barh

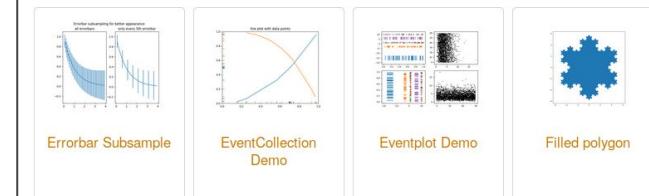


Plotting categorical variables

Plotting the coherence of two signals

CSD Demo

Errorbar limit selection



Errorbar Subsample

EventCollection Demo

Eventplot Demo

Filled polygon

To install this package with conda run one of the following:

```
conda install -c conda-forge matplotlib
```

<https://matplotlib.org/gallery/index.html>

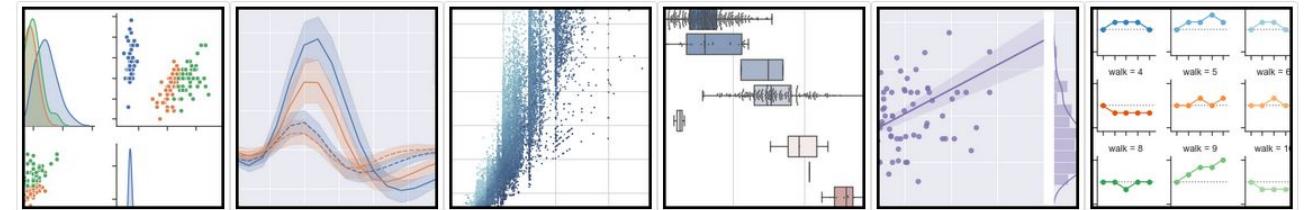
Visualisation 3 / 10



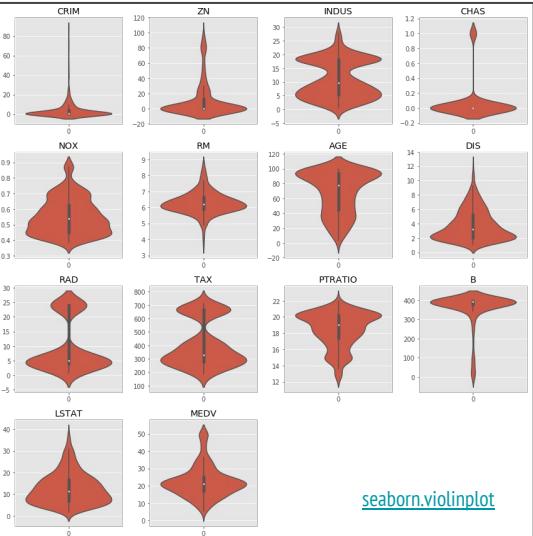
seaborn

seaborn: statistical data visualization

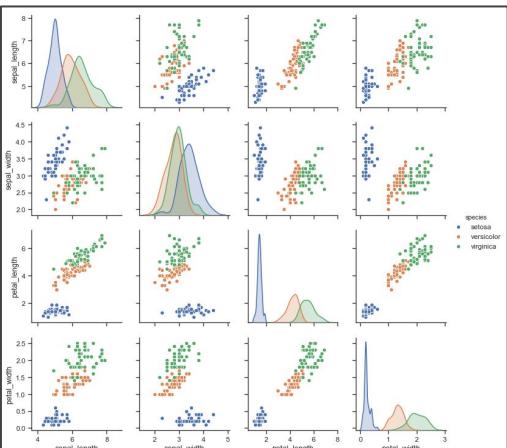
<https://seaborn.pydata.org/>



Plus facile à utiliser, avec peu de paramètre à fournir et fait pour être utilisé par des statisticiens (et data scientist). Les graphiques sont souvent plus "jolis"/"Pro".



[seaborn.violinplot](#)



[seaborn.pairplot](#)

To install this package with conda run:

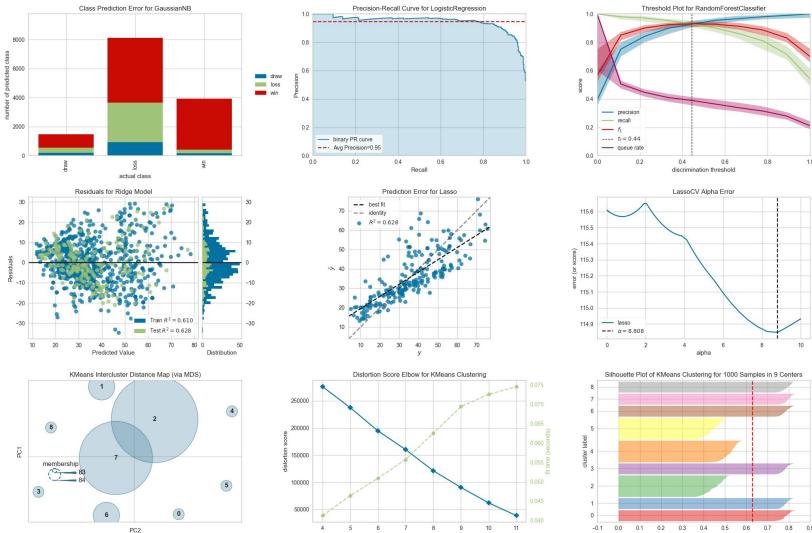
```
conda install -c anaconda seaborn
```

Visualisation 4 / 10

Yellowbrick

build passing build passing coverage 90% lgtn alerts 13 code quality: python A pypi package 1.2 docs passing
DOI 10.5281/zenodo.1206239 JOSS 10.21105/oss.01075 launch binder

Visual analysis and diagnostic tools to facilitate machine learning model selection.



Feature Visualization

- :doc:`api/features/rankd` : pairwise ranking of features to detect relationships
- :doc:`api/features/pcoords` : horizontal visualization of instances
- :doc:`Radial Visualization <api/features/radviz>` : separation of instances around a circular plot
- :doc:`api/features/pca` : projection of instances based on principal components
- :doc:`api/features/manifold` : high dimensional visualization with manifold learning
- :doc:`Joint Plots <api/features/jointplot>` : direct data visualization with feature selection

Classification Visualization

- :doc:`api/classifier/class_prediction_error` : shows error and support in classification
- :doc:`api/classifier/classification_report` : visual representation of precision, recall, and F1
- :doc:`ROC/AUC Curves <api/classifier/roc_auc>` : receiver operator characteristics and area under the curve
- :doc:`api/classifier/prcurve` : precision vs recall for different probability thresholds
- :doc:`Confusion Matrices <api/classifier/confusion_matrix>` : visual description of class decision making
- :doc:`Discrimination Threshold <api/classifier/threshold>` : find a threshold that best separates binary classes

Regression Visualization

- :doc:`api/regressor/peplot` : find model breakdowns along the domain of the target
- :doc:`api/regressor/residuals` : show the difference in residuals of training and test data
- :doc:`api/regressor/alphas` : show how the choice of alpha influences regularization
- :doc:`api/regressor/influence` : show the influence of instances on linear regression

Clustering Visualization

- :doc:`K-Elbow Plot <api/cluster/elbow>` : select k using the elbow method and various metrics
- :doc:`Silhouette Plot <api/cluster/silhouette>` : select k by visualizing silhouette coefficient values
- :doc:`api/cluster/lcdm` : show relative distance and size/importance of clusters

Model Selection Visualization

- :doc:`api/model_selection/validation_curve` : tune a model with respect to a single hyperparameter
- :doc:`api/model_selection/learning_curve` : show if a model might benefit from more data or less complexity
- :doc:`api/model_selection/importances` : rank features by importance or linear coefficients for a specific model
- :doc:`api/model_selection/rfecv` : find the best subset of features based on importance

Target Visualization

- :doc:`api/target/binning` : generate a histogram with vertical lines showing the recommended value point to bin the data into evenly distributed bins
- :doc:`api/target/class_balance` : see how the distribution of classes affects the model
- :doc:`api/target/feature_correlation` : display the correlation between features and dependent variables

Text Visualization

- :doc:`Term Frequency <api/text/freqdist>` : visualize the frequency distribution of terms in the corpus
- :doc:`api/text/tse` : use stochastic neighbor embedding to project documents
- :doc:`api/text/dispersion` : visualize how key terms are dispersed throughout a corpus
- :doc:`api/text/umap_vis` : plot similar documents closer together to discover clusters
- :doc:`api/text/postag` : plot the counts of different parts-of-speech throughout a tagged corpus

Visualisation 5 / 10

dtreeviz : Decision Tree Visualization

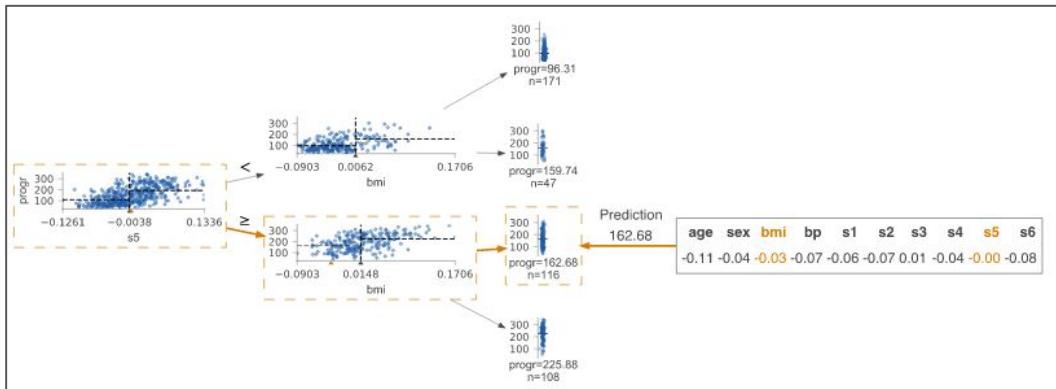
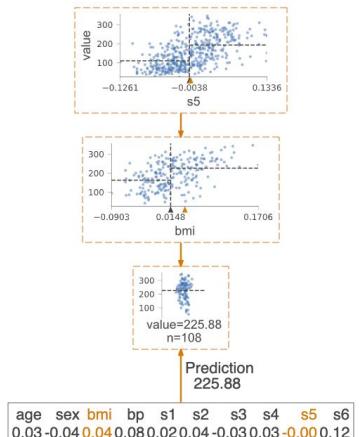
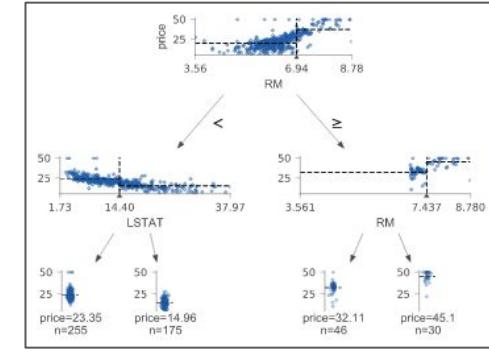
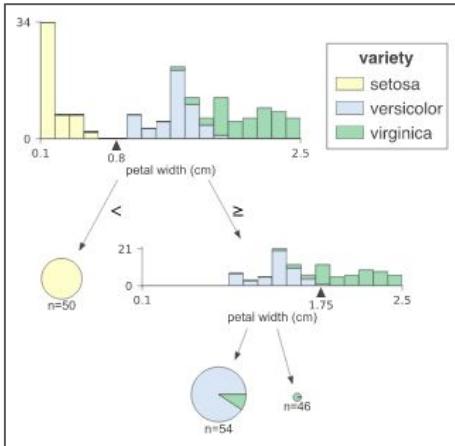
Description

A python library for decision tree visualization and model interpretation. Currently supports scikit-learn, XGBoost and Spark MLlib trees.

Authors:

- Terence Parr, a professor in the University of San Francisco's data science program
- Tudor Lapusan
- Prince Grover

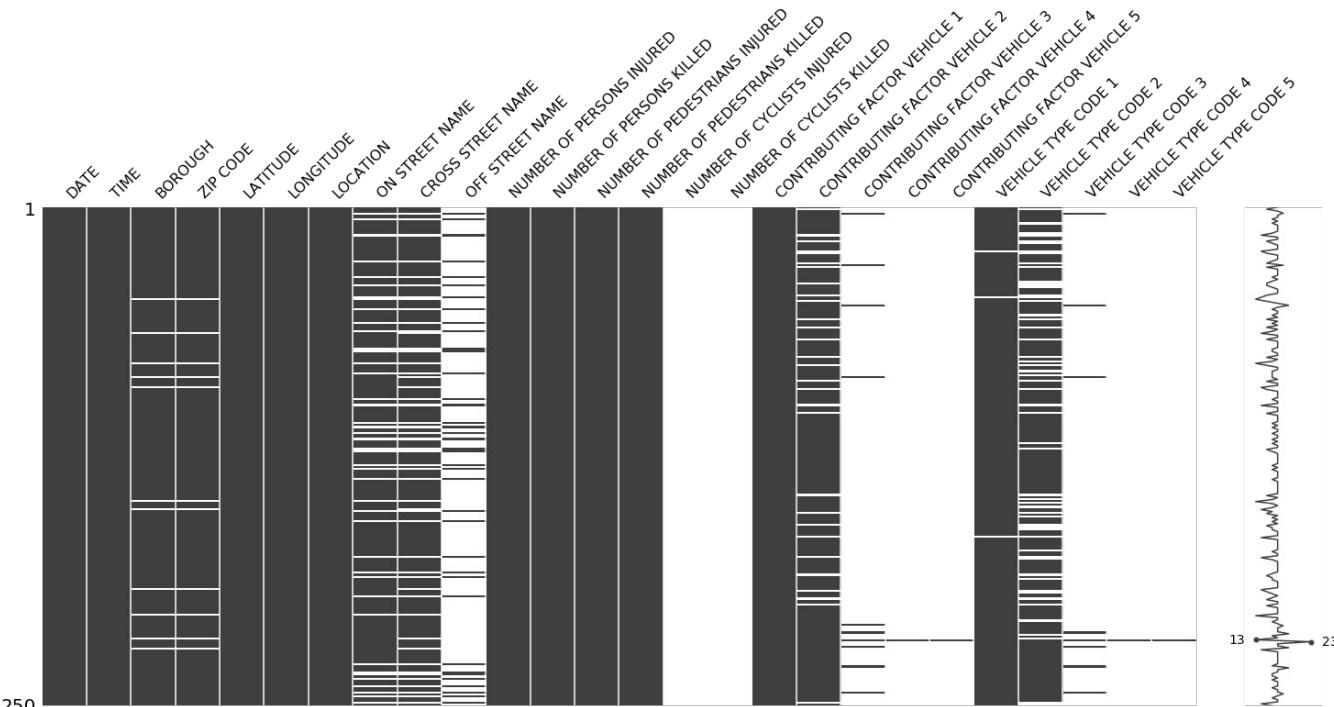
See [How to visualize decision trees](#) for deeper discussion of our decision tree visualization library and the visual design decisions we made.



<https://github.com/parrt/dtreeviz>

Visualisation 6 / 10

Visualisation des données absentes.



missingno pypi v0.4.2 python 3.4+ status stable license MIT doi 10.21105/joss.00547+

Messy datasets? Missing values? `missingno` provides a small toolset of flexible and easy-to-use missing data visualizations and utilities that allows you to get a quick visual summary of the completeness (or lack thereof) of your dataset. Just `pip install missingno` to get started.

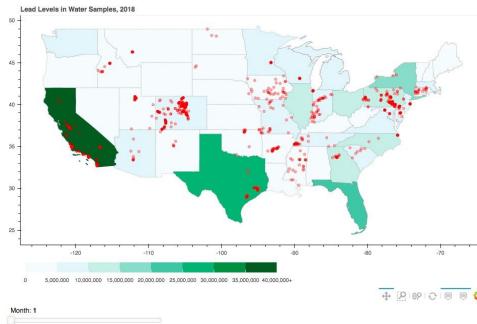
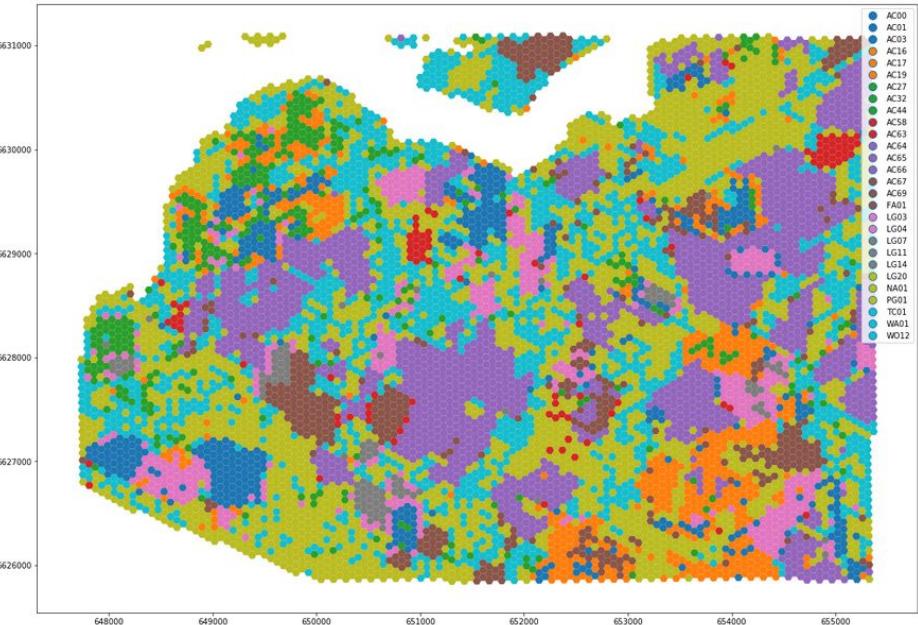
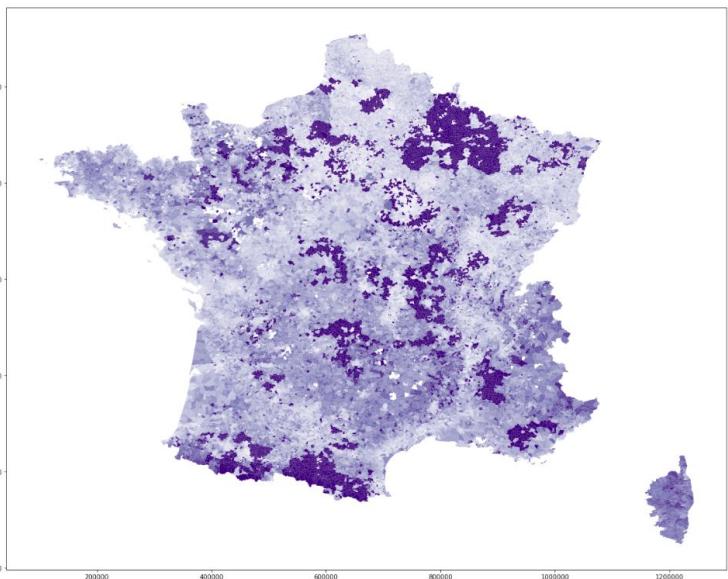
Visualisation 7 / 10

Geopandas

Project description

GeoPandas is a project to add support for geographic data to [pandas](#) objects.

The goal of GeoPandas is to make working with geospatial data in python easier. It combines the capabilities of [pandas](#) and [shapely](#), providing geospatial operations in pandas and a high-level interface to multiple geometries to shapely. GeoPandas enables you to easily do operations in python that would otherwise require a spatial database such as PostGIS.



Visualisation sur de cartes
(formes)

<https://geopandas.org/>

Visualisation 8 / 10

pypi v2.0.0 build passing license MIT

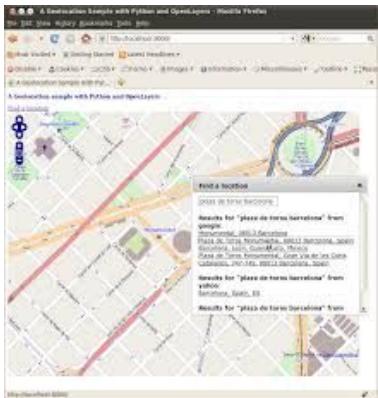
geopy is a Python client for several popular geocoding web services.

geopy makes it easy for Python developers to locate the coordinates of addresses, cities, countries, and landmarks across the globe using third-party geocoders and other data sources.

geopy includes geocoder classes for the OpenStreetMap Nominatim, Google Geocoding API (V3), and many other geocoding services. The full list is available on the [Geocoders doc section](#). Geocoder classes are located in `geopy.geocoders`.

geopy is tested against CPython (versions 3.5, 3.6, 3.7, 3.8) and PyPy3. geopy 1.x line also supported CPython 2.7, 3.4 and PyPy2.

© geopy contributors 2006-2018 (see AUTHORS) under the [MIT License](#).



geopy is a Python client for several popular geocoding web services.

geopy makes it easy for Python developers to locate the coordinates of addresses, cities, countries, and landmarks across the globe using third-party geocoders and other data sources.

geopy includes geocoder classes for the OpenStreetMap Nominatim, Google Geocoding API (V3), and many other geocoding services. The full list is available on the [Geocoders doc section](#). Geocoder classes are located in `geopy.geocoders`.

geopy is tested against CPython (versions 3.5, 3.6, 3.7, 3.8) and PyPy3. geopy 1.x line also supported CPython 2.7, 3.4 and PyPy2.

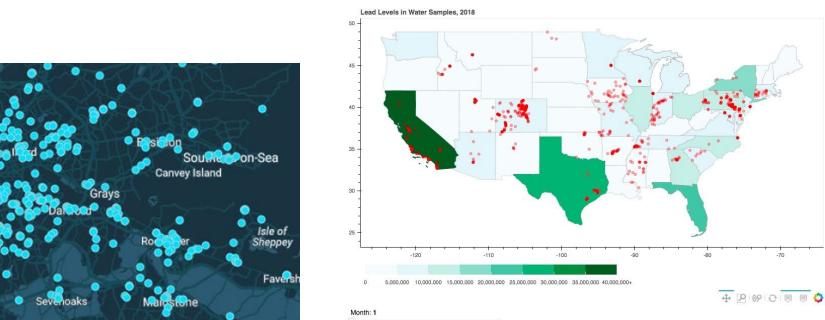
© geopy contributors 2006-2018 (see AUTHORS) under the [MIT License](#).



OpenStreetMap

in Python

Set Longitudes and Latitudes



OpenStreetMap

downloads 5k

`openstreetmap` is a pure Python library that provides an easy way to extracting OpenStreetMap coordinates by name or relation id.

<https://github.com/Mywayking/openstreetmap>

<https://github.com/geopy/geopy>

Visualisation 9 / 10

Commande dans Notebook

```
%matplotlib notebook
```

ipywidgets: Interactive HTML Widgets

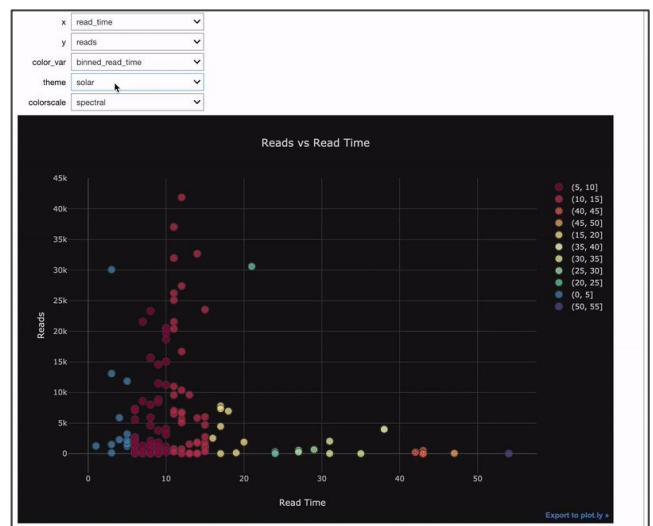
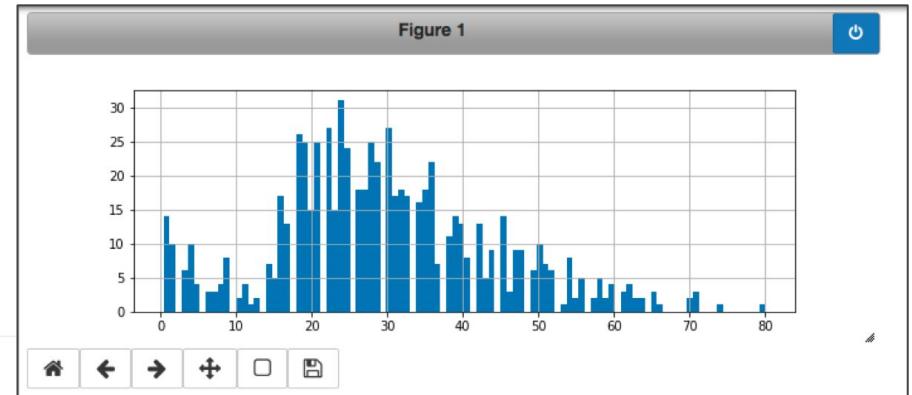
pypi v7.5.1 build passing docs passing gitter join chat launch binder

ipywidgets are [Interactive HTML widgets](#) for Jupyter notebooks and the IPython kernel.

Notebooks come alive when interactive widgets are used. Users gain control of their data and can visualize changes in the data.

IPywidgets <https://matplotlib.org/>

Création de Widgets interactif dans l'output de NoteBook. Cela permet de pouvoir d'explorer plus facilement / rapidement avec les données.



Visualisation 10 / 10

User Showcase

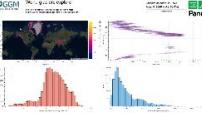
Dask



Dask is a tool for scaling out PyData projects like NumPy, Pandas, Scikit-Learn, and RAPIDS. It is supported by Nvidia, Quansight, and Anaconda.

The [Dask Dashboard](#) is a diagnostic tool that helps you monitor and debug live cluster performance.

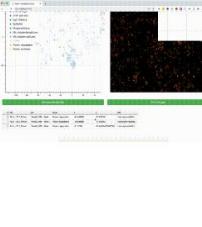
Panel



Panel is a tool for polished data presentation that utilizes the Bokeh server. It is created and supported by Anaconda.

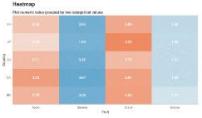
Panel makes it simple to create custom interactive web apps and dashboards by connecting user-defined widgets to plots, images, tables, or text.

Microscopium

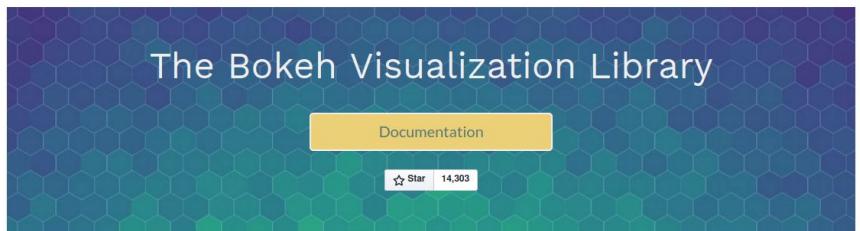


Microscopium is a project maintained by researchers at Monash University. It allows researchers to discover new gene or drug functions by exploring large image datasets with Bokeh's interactive tools.

Chartify



Chartify is an opinionated high-level charting API built on top of Bokeh, created by Spotify. With smart default styles, consistent tidy data format, and a simple API, it's easy for you to concentrate on your work.



The Bokeh Visualization Library

Documentation

Star 14,303

Bokeh at a Glance

Flexible

Bokeh makes it simple to create common plots, but also can handle custom or specialized use-cases.

Interactive

Tools and widgets let you and your audience probe "what if" scenarios or drill-down into the details of your data.

Productive

Work in Python close to all the PyData tools you are already familiar with.

Powerful

You can always add custom JavaScript to support advanced or specialized cases.

Shareable

Plots, dashboards, and apps can be published in web pages or Jupyter notebooks.

Open Source

Everything, including the Bokeh server, is BSD licensed and available on [GitHub](#).

<https://bokeh.org/>

Jeux de données



Jeux de Données 1 /



WIKIPEDIA
The Free Encyclopedia

kaggle

<https://www.kaggle.com/datasets>

List of datasets for machine-learning research

https://en.wikipedia.org/wiki/List_of_datasets_for_machine-learning_research

Mon  coach Data

[https://moncoachdata.com/blog/datasets-
projet-data-science/](https://moncoachdata.com/blog/datasets-projet-data-science/)



Machine Learning Repository

<https://archive.ics.uci.edu/ml/>

Awesome Public Datasets Collection

<https://github.com/awesomedata/awesome-public-datasets>



Jeux de Données 2 / 2 : Officiels



EU **Open Data** Portal

<https://data.europa.eu/euodp/data/dataset>



<https://www.data.gov/>



<https://data.gov.in/>



data.gouv.fr

<https://www.data.gouv.fr/fr/>



<https://data.gov.hk/en/>



<https://dataportal.opendataforafrica.org/>



<https://datos.gob.cl/dataset>



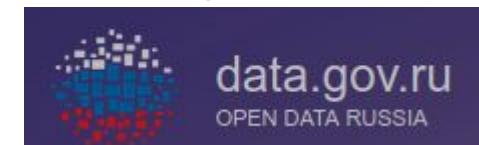
stats sa

Department:
Statistics South Africa
REPUBLIC OF SOUTH AFRICA

<http://www.statssa.gov.za/>



<https://data.strasbourg.eu/pages/accueil/>



<https://data.gov.ru/>



<https://ec.europa.eu/eurostat/data/database>



Manipulation des données

<https://www.datacamp.com/community/data-science-cheatsheets>



Python For Data Science Cheat Sheet

Pandas Basics

Learn Python for Data Science Interactively at www.DataCamp.com



Pandas

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

pandas



Use the following import convention:

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

Pandas Data Structures

Series

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

a	3
b	-5
c	7
d	4

Index

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

DataFrame

Columns

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

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],
   'Capital': ['Brussels', 'New Delhi', 'Brasilia'],
   'Population': [11190846, 1303171035, 207847528]}
```

```
>>> df = pd.DataFrame(data,
   columns=['Country', 'Capital', 'Population'])
```

I/O

Read and Write to CSV

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

Read and Write to Excel

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

Asking For Help

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

Selection

Also see NumPy Arrays

Getting

```
>>> s['b']
```

-5

```
>>> df[1:]
```

	Country	Capital	Population
1	India	New Delhi	1303171035
2	Brazil	Brasilia	207847528

Get one element

Get subset of a DataFrame

Selecting, Boolean Indexing & Setting

By Position

```
>>> df.iloc[[0], [0]]
```

'Belgium'

```
>>> df.iat[[0], [0]]
```

'Belgium'

By Label

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

'Belgium'

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

'Belgium'

By Label/Position

```
>>> df.ix[2]
```

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

```
>>> df.ix[:, 'Capital']
```

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

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

'New Delhi'

Boolean Indexing

```
>>> s[~(s > 1)]
```

```
>>> s[(s < -1) | (s > 2)]
```

```
>>> df[df['Population']>1200000000]
```

Setting

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

Select single value by row & column

Select single value by row & column labels

Select single row of subset of rows

Select a single column of subset of columns

Select rows and columns

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

Set index a of Series s to 6

Dropping

```
>>> s.drop(['a', 'c'])
```

Drop values from rows (axis=0)

```
>>> df.drop('Country', axis=1)
```

Drop values from columns (axis=1)

Sort & Rank

```
>>> df.sort_index()
```

Sort by labels along an axis

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

Sort by the values along an axis

```
>>> df.rank()
```

Assign ranks to entries

Retrieving Series/DataFrame Information

Basic Information

```
>>> df.shape
```

(rows,columns)

```
>>> df.index
```

Describe index

```
>>> df.columns
```

Describe DataFrame columns

```
>>> df.info()
```

Info on DataFrame

```
>>> df.count()
```

Number of non-NA values

Summary

```
>>> df.sum()
```

Sum of values

```
>>> df.cumsum()
```

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
```

Apply function

```
>>> df.apply(f)
```

Apply function element-wise

Data Alignment

Internal Data Alignment

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

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
```

```
>>> s + s3
```

a 10.0

b NaN

c 5.0

d 7.0

Arithmetic Operations with Fill Methods

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

```
>>> s.add(s3, fill_value=0)
```

a 10.0

b -5.0

c 5.0

d 7.0

```
>>> s.sub(s3, fill_value=2)
```

```
>>> s.div(s3, fill_value=4)
```

```
>>> s.mul(s3, fill_value=3)
```



Python For Data Science Cheat Sheet

Pandas

Learn Python for Data Science Interactively at www.DataCamp.com



Reshaping Data

Pivot

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

Spread rows into columns

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

Where

```
>>> s.where(s > 0)
Query
>>> df6.query('second > first')
```

Setting/Resetting Index

```
>>> df.set_index('Country')
>>> df4 = df.reset_index()
>>> df = df.rename(index=str,
                  columns={"Country": "cntry",
                            "Capital": "cptl",
                            "Population": "popln"})
```

Set the index
Reset the index
Rename DataFrame

Pivot Table

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

Spread rows into columns

Stack / Unstack

```
>>> stacked = df5.stack()
>>> stacked.unstack()
```

Pivot a level of column labels
Pivot a level of index labels

	0	1
1	0.233482	0.390959
2	0.184713	0.237102
3	0.433522	0.429401

	0	1
1	0.233482	0.390959
2	0.184713	0.237102
3	0.433522	0.429401

Unstacked Stacked

Melt

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

Gather columns into rows

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

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

Iteration

```
>>> df.iteritems()
>>> df.iterrows()
```

(Column-index, Series) pairs
(Row-index, Series) pairs

Advanced Indexing

Selecting

```
>>> df3.loc[:, (df3>1).any()]
>>> df3.loc[:,(df3>1).all()]
>>> df3.loc[:,df3.isnull().any()]
>>> df3.loc[:,df3.notnull().all()]
Indexing With isin
>>> df[(df.Country.isin(df2.Type))]
>>> df3.filter(items="a","b")
>>> df.select(lambda x: not x%5)
```

Select cols with any vals >1
Select cols with vals > 1
Select cols with NaN
Select cols without NaN

Find same elements
Filter on values
Select specific elements

Where

```
>>> s.where(s > 0)
Query
>>> df6.query('second > first')
```

Reindexing

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

Forward Filling

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

Country	Capital	Population
0	Belgium	11190846
1	India	1303171035
2	Brazil	207847528
3	Brazil	207847528

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

	0	3
0	1	3
1	1	3
2	1	3
3	1	3
4	1	3

MultIndexing

```
>>> 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)
>>> df6.set_index(["Date", "Type"])
```

Duplicate Data

```
>>> s3.unique()
>>> df2.duplicated('Type')
>>> df2.drop_duplicates('Type', keep='last')
>>> df.index.duplicated()
```

Return unique values
Check duplicates
Drop duplicates
Check index duplicates

Grouping Data

Aggregation

```
>>> df2.groupby(by=['Date', 'Type']).mean()
>>> df4.groupby(level=0).sum()
>>> df4.groupby(level=0).agg([('a':lambda x:sum(x)/len(x),
                                'b': np.sum)})
```

Transformation

```
>>> customSum = lambda x: (x*x2)
>>> df4.groupby(level=0).transform(customSum)
```

Missing Data

```
>>> df.dropna()
>>> df3.fillna(df3.mean())
>>> df2.replace("a", "f")
```

Drop Na values
Fill Na values with a predetermined value
Replace values with others

Also see NumPy Arrays

Select cols with any vals >1
Select cols with vals > 1
Select cols with NaN
Select cols without NaN

Find same elements
Filter on values
Select specific elements

Subset the data

Query DataFrame

Combining Data

data1	X1	X2
a	11.432	20.784
b	1.303	NaN
c	99.906	20.784

data2	X1	X3
a	11.432	20.784
b	NaN	NaN
d	NaN	20.784

Merge

```
>>> pd.merge(data1,
             data2,
             how='left',
             on='X1')
```

X1	X2	X3
a	11.432	20.784
b	1.303	NaN
c	99.906	NaN

X1	X2	X3
a	11.432	20.784
b	1.303	NaN
d	NaN	20.784

X1	X2	X3
a	11.432	20.784
b	1.303	NaN
c	99.906	NaN
d	NaN	20.784

Join

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

Concatenate

Vertical

```
>>> s.append(s2)
```

Horizontal/Vertical

```
>>> pd.concat([s,s2],axis=1, keys=['One','Two'])
>>> pd.concat([data1, data2], axis=1, join='inner')
```

Dates

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

Visualization

Also see Matplotlib

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



DataCamp

Learn Python for Data Science Interactively



Data Wrangling

with pandas
Cheat Sheet

<http://pandas.pydata.org>

Syntax – Creating DataFrames

	a	b	c
1	4	7	10
2	5	8	11
3	6	9	12

```
df = pd.DataFrame(
    {"a" : [4, 5, 6],
     "b" : [7, 8, 9],
     "c" : [10, 11, 12]},
    index = [1, 2, 3])
Specify values for each column.
```

```
df = pd.DataFrame(
    [[4, 7, 10],
     [5, 8, 11],
     [6, 9, 12]],
    index=[1, 2, 3],
    columns=['a', 'b', 'c'])
Specify values for each row.
```

n	v	a	b	c
d	1	4	7	10
e	2	5	8	11

```
df = pd.DataFrame(
    {"a" : [4, 5, 6],
     "b" : [7, 8, 9],
     "c" : [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [('d',1),('d',2),('e',2)],
        names=['n', 'v']))
Create DataFrame with a MultiIndex
```

Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```
df = (pd.melt(df)
      .rename(columns={'variable' : 'var',
                      'value' : 'val'})
      .query('val >= 200'))
```

Tidy Data – A foundation for wrangling in pandas

In a tidy data set:

Each variable is saved in its own column

Each observation is saved in its own row

Tidy data complements pandas's **vectorized operations**. pandas will automatically preserve observations as you manipulate variables. No other format works as intuitively with pandas.

M * A

Reshaping Data – Change the layout of a data set

 pd.melt(df) Gather columns into rows.	 df.pivot(columns='var', values='val') Spread rows into columns.
 pd.concat([df1, df2]) Append rows of DataFrames	 pd.concat([df1, df2], axis=1) Append columns of DataFrames

df.sort_values('mpg')
Order rows by values of a column (low to high).

df.sort_values('mpg', ascending=False)
Order rows by values of a column (high to low).

df.rename(columns = {'y': 'year'})
Rename the columns of a DataFrame

df.sort_index()
Sort the index of a DataFrame

df.reset_index()
Reset index of DataFrame to row numbers, moving index to columns.

df.drop(columns=['Length', 'Height'])
Drop columns from DataFrame

Subset Observations (Rows)



```
df[df.Length > 7]
Extract rows that meet logical criteria.

df.drop_duplicates()
Remove duplicate rows (only considers columns).

df.head(n)
Select first n rows.

df.tail(n)
Select last n rows.
```

```
df.sample(frac=0.5)
Randomly select fraction of rows.

df.sample(n=10)
Randomly select n rows.

df.iloc[10:20]
Select rows by position.

df.nlargest(n, 'value')
Select and order top n entries.

df.nsmallest(n, 'value')
Select and order bottom n entries.
```

Subset Variables (Columns)



```
df[['width', 'length', 'species']]
Select multiple columns with specific names.

df['width'] or df.width
Select single column with specific name.

df.filter(regex='regex')
Select columns whose name matches regular expression regex.
```

regex (Regular Expressions) Examples

'.'	Matches strings containing a period '.'.
'Length\$'	Matches strings ending with word 'Length'
'Sepal'	Matches strings beginning with the word 'Sepal'
'x{1-5}\$'	Matches strings beginning with 'x' and ending with 1,2,3,4,5
'^^(?!Species).*\$'	Matches strings except the string 'Species'

```
df.loc[:, 'x2':'x4']
Select all columns between x2 and x4 (inclusive).

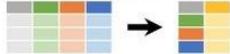
df.iloc[:, 1,2,5]
Select columns in positions 1, 2 and 5 (first column is 0).

df.loc[df['a'] > 10, ['a', 'c']]
Select rows meeting logical condition, and only the specific columns.
```

Logic in Python (and pandas)		
<	Less than	!=
>	Greater than	df.column.isin(values)
==	Equals	pd.isnull(obj)
<=	Less than or equals	pd.notnull(obj)
>=	Greater than or equals	&, , ~, ^, df.any(), df.all()
		Not equal to
		Group membership
		Is NaN
		Is not NaN
		Logical and, or, not, xor, any, all

Summarize Data

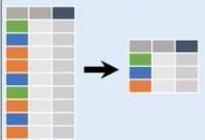
```
df['Length'].value_counts()
    Count number of rows with each unique value of variable
len(df)
    # of rows in DataFrame.
len(df['w'].unique())
    # of distinct values in a column.
df.describe()
    Basic descriptive statistics for each column (or GroupBy)
```



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

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

Group Data



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

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

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

size()	Size of each group.
agg(function)	Aggregate group using function.

Windows

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

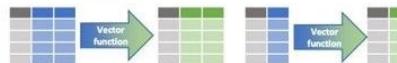
Handling Missing Data

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

Make New Variables



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



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

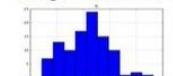
max(axis=1)	Element-wise max.
clip(lower=-10,upper=10)	Element-wise min.
abs()	Absolute value.
Trim values at input thresholds	

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

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

Plotting

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



Combine Data Sets

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

Standard Joins

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

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

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

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

Filtering Joins

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

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

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

Set-like Operations

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

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

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

Python For Data Science Cheat Sheet

Scikit-Learn

Learn Python for data science interactively at www.DataCamp.com



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[:, :-1], 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','F','M','M','F','F'])
>>> X[X < 0.7] = 0
```

Training And Test Data

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

Preprocessing The Data

Standardization

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

Normalization

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

Binarization

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

Create Your Model

Supervised Learning Estimators

Linear Regression

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

Support Vector Machines (SVM)

```
>>> from sklearn.svm import SVC
>>> svc = SVC(kernel='linear')
```

Naive Bayes

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

KNN

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

Unsupervised Learning Estimators

Principal Component Analysis (PCA)

```
>>> from sklearn.decomposition import PCA
>>> pca = PCA(n_components=0.95)
```

K-Means

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

Model Fitting

Supervised learning

```
>>> lr.fit(X, y)
>>> knn.fit(X_train, y_train)
>>> svc.fit(X_train, y_train)
```

Unsupervised Learning

```
>>> k_means.fit(X_train)
>>> pca_model = pca.fit_transform(X_train)
```

Fit the model to the data

Fit the model to the data

Fit to data, then transform it

Prediction

Supervised Estimators

```
>>> y_pred = svc.predict(np.random.random(2,5))
>>> y_pred = lr.predict(X_test)
>>> y_pred = knn.predict_proba(X_test)
```

Unsupervised Estimators

```
>>> y_pred = k_means.predict(X_test)
```

Predict labels

Predict labels

Estimate probability of a label

Predict labels in clustering algos

Evaluate Your Model's Performance

Classification Metrics

Accuracy Score

```
>>> knn.score(X_test, y_test)
>>> from sklearn.metrics import accuracy_score
>>> accuracy_score(y_test, y_pred)
```

Estimator score method

Metric scoring functions

Classification Report

```
>>> from sklearn.metrics import classification_report
>>> print(classification_report(y_test, y_pred))
```

Confusion Matrix

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

Precision, recall, f1-score and support

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)
```

R² 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=knk,
...                               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 Cheat Sheet

Keras

Learn Python for data science interactively at www.DataCamp.com



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

Data

Also see NumPy, Pandas & Scikit-Learn

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

Keras Data Sets

```
>>> 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-diabetes/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'))
```

Also see NumPy & Scikit-Learn

Train and Test Sets

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

Standardization/Normalization

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

Inspect Model

```
>>> model.output_shape
>>> model.summary()
>>> model.get_config()
>>> model.get_weights()
```

Model output shape
Model summary representation
Model configuration
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_train4,
                y_train4,
                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)
```

Prediction

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

Importing Data

Learn Python for data science interactively at www.DataCamp.com



Importing Data in Python

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

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

Help

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

Text Files

Plain Text Files

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

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

Using the context manager with

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

Read a single line

Table Data: Flat Files

Importing Flat Files with numpy

Files with one data type

```
>>> filename = 'mnist.txt'  
>>> data = np.loadtxt(filename,  
                     delimiter=',',  
                     skiprows=2,  
                     usecols=[0, 21],  
                     dtype=str)
```

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

Files with mixed data types

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

Look for column header

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

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

Importing Flat Files with pandas

```
>>> filename = 'winequality-red.csv'  
>>> data = pd.read_csv(filename,  
                     nrows=5,  
                     header=None,  
                     sep=';'  
                     comment="#",  
                     na_values=[""])
```

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

Excel Spreadsheets

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

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

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

SAS Files

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

Stata Files

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

Relational Databases

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

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

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

Querying Relational Databases

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

Using the context manager with

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

Querying relational databases with pandas

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

Exploring Your Data

NumPy Arrays

<code>>>> data_array.dtype</code>	Data type of array elements
<code>>>> data_array.shape</code>	Array dimensions
<code>>>> len(data_array)</code>	Length of array

pandas DataFrames

<code>>>> df.head()</code>	Return first DataFrame rows
<code>>>> df.tail()</code>	Return last DataFrame rows
<code>>>> df.index</code>	Describe index
<code>>>> df.columns</code>	Describe DataFrame columns
<code>>>> df.info()</code>	Info on DataFrame
<code>>>> data_array = data.values</code>	Convert a DataFrame to an NumPy array

Pickled Files

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

HDF5 Files

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

Matlab Files

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

Exploring Dictionaries

Accessing Elements with Functions

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

Accessing Data Items with Keys

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

Navigating Your FileSystem

Magic Commands

```
:ls  
%cd ..  
$pwd
```

List directory contents of files and directories
Change current working directory
Return the current working directory path

os Library

```
>>> import os  
>>> path = "/usr/tmp"  
>>> wd = os.getcwd()  
>>> os.listdir(wd)  
>>> os.chdir(path)  
>>> os.rename("test1.txt", "test2.txt")  
>>> os.remove("test1.txt")  
>>> os.mkdir("newdir")
```

Store the name of current directory in a string
Output contents of the directory in a list
Change current working directory
Rename a file
Delete an existing file
Create a new directory



Python For Data Science Cheat Sheet

SciPy - Linear Algebra

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



SciPy

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



Interacting With NumPy

[Also see NumPy](#)

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

Index Tricks

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

Shape Manipulation

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

Polynomials

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

Vectorizing Functions

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

Type Handling

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

Linear Algebra

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

[Also see NumPy](#)

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  
>>> linalg.inv(A)  
>>> A.T  
>>> A.H  
>>> np.trace(A)
```

Norm

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

Rank

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

Determinant

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

Solving linear problems

```
>>> linalg.solve(A,b)  
>>> E = np.mat(a).T  
>>> linalg.lstsq(D,E)
```

Generalized inverse

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

Inverse

```
Inverse  
Transpose matrix  
Conjugate transposition  
Trace
```

Frobenius norm

```
L1 norm (max column sum)  
L inf norm (max row sum)
```

Matrix rank

Determinant

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

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

Creating Sparse Matrices

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

Sparse Matrix Routines

Inverse

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

Norm

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

Solving linear problems

```
>>> sparse.linalg.spsolve(H,I) Solver for sparse matrices
```

Inverse

```
Inverse
```

Norm

```
Norm
```

```
Solver for sparse matrices
```

Sparse Matrix Functions

```
>>> sparse.linalg.expm(I) Sparse matrix exponential
```

Asking For Help

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

Matrix Functions

Addition

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

Subtraction

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

Division

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

Multiplication

```
>>> np.multiply(D,A)  
>>> np.dot(A,D)  
>>> np.vdot(A,D)  
>>> np.inner(A,D)  
>>> np.outer(A,D)  
>>> np.tensordot(A,D)  
>>> np.kron(A,D)
```

Exponential Functions

```
>>> linalg.expm(A)  
>>> linalg.expm2(A)  
>>> linalg.expm3(D)
```

Logarithm Function

```
>>> linalg.logm(A)
```

Trigonometric Functions

```
>>> linalg.sinm(D)
```

```
>>> linalg.cosm(D)
```

```
>>> linalg.tanm(A)
```

Hyperbolic Trigonometric Functions

```
>>> linalg.sinh(D)
```

```
>>> linalg.cosh(D)
```

```
>>> linalg.tanh(A)
```

Matrix Sign Function

```
>>> np.sign(A)
```

Matrix Square Root

```
>>> linalg.sqrtm(A)
```

Arbitrary Functions

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

Evaluate matrix function

Decompositions

Eigenvalues and Eigenvectors

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

```
>>> 11, 12 = la
```

```
>>> v[:,0]
```

```
>>> v[:,1]
```

```
>>> linalg.eigvals(A)
```

Singular Value Decomposition

```
>>> U,s,Vh = linalg.svd(B)
```

```
>>> M,N = B.shape
```

```
>>> Sig = linalg.diagsvd(s,M,N)
```

LU Decomposition

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

LU Decomposition

Sparse Matrix Decompositions

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

Eigenvalues and eigenvectors

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

SVD



Python For Data Science Cheat Sheet

NumPy Basics

Learn Python for Data Science Interactively at www.DataCamp.com



NumPy

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

Use the following import convention:

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



NumPy Arrays

1D array

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

2D array

1.5	2	3
4	5	6

3D array

axis 0	axis 1	axis 2

Creating Arrays

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

Initial Placeholders

```
>>> np.zeros((3,4))
>>> np.ones((2,3,4),dtype=np.int16)
>>> d = np.arange(10,25)
>>> np.linspace(0,2,9)
>>> e = np.full((2,2),7)
>>> f = np.eye(2)
>>> np.random.random((2,2))
>>> np.empty((3,2))
```

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

I/O

Saving & Loading On Disk

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

Saving & Loading Text Files

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

Data Types

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

Inspecting Your Array

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

Asking For Help

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

Array Mathematics

Arithmetic Operations

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

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

Comparison

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

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

Aggregate Functions

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

Copying Arrays

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

Sorting Arrays

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

Subsetting, Slicing, Indexing

Also see [Lists](#)

Subsetting

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



Select the element at the 2nd index
Select the element at row 1 column 2 (equivalent to `b[1][2]`)

Slicing

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



Select items at index 0 and 1
Select items at rows 0 and 1 in column 1
Select all items at row 0 (equivalent to `b[0,:,:]`)
Same as `[1,:,:]`

Reversed array `a`

Select elements from `a` less than 2

```
>>> a[a<2]
array([1])
>>> b[[1, 0, 1, 0], [0, 1, 2, 0]]
array([[ 1.,  2.,  3.,  4.],
      [ 1.5,  2.,  3.,  4.5],
      [ 1.5,  2.,  3.,  4.5],
      [ 1.5,  2.,  3.,  4.5]])
```

Select elements `(1,0),(0,1),(1,2)` and `(0,0)`
Select a subset of the matrix's rows and columns

Array Manipulation

Transposing Array

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

Permute array dimensions
Permute array dimensions

Changing Array Shape

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

Flatten the array
Reshape, but don't change data

Adding/Removing Elements

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

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

Combining Arrays

```
>>> np.concatenate((a,d),axis=0)
array([[ 1.,  2.,  3., 10., 15., 20.])
>>> np.vstack((a,b))
array([[ 1.,  2.,  3.,  4.],
      [ 1.5,  2.,  3.,  4.5],
      [ 1.5,  2.,  3.,  4.5],
      [ 1.5,  2.,  3.,  4.5]])
```

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

Create stacked column-wise arrays

Create stacked column-wise arrays

Split the array horizontally at the 3rd index

Split the array vertically at the 2nd index



Python For Data Science Cheat Sheet

3) Renderers & Visual Customizations

Bokeh

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



Plotting With Bokeh

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



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



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

1. Prepare some data:

Python lists, NumPy arrays, Pandas DataFrames and other sequences of values

2. Create a new plot

3. Add renderers for your data, with visual customizations
4. Specify where to generate the output
5. Show or save the results

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

Glyphs

Scatter Markers

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

Line Glyphs

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

Customized Glyphs

Also see Data

```
>>> 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')
```

Legend Location

Inside Plot Area

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

Outside Plot Area

```
>>> from bokeh.models import Legend
>>> r1 = p2.pasterine(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)
```

4 Output & Export

Notebook

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

HTML

Standalone HTML

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

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

Components

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

PNG

```
>>> from bokeh.io import export_png
>>> export_png(p, filename="plot.png")
```

SVG

```
>>> from bokeh.io import export_svgs
>>> p.output_backend = "svg"
>>> export_svgs(p, filename="plot.svg")
```

5 Show or Save Your Plots

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



Python For Data Science Cheat Sheet

Seaborn

Learn Data Science Interactively at www.DataCamp.com



Statistical Data Visualization With Seaborn

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

Make use of the following aliases to import the libraries:

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

The basic steps to creating plots with Seaborn are:

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

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

Also see Lists, NumPy & Pandas

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

Seaborn also offers built-in data sets:

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

2 Figure Aesthetics

```
>>> f, ax = plt.subplots(figsize=(5,6))
```

Create a figure and one subplot

Seaborn styles

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

(Re)set the seaborn default
Set the matplotlib parameters
Set the matplotlib parameters
Return a dict of params or use with
with to temporarily set the style

3 Plotting With Seaborn

Axis Grids

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

Subplot grid for plotting conditional relationships

Draw a categorical plot onto a Facetgrid

Plot data and regression model fits across a FacetGrid

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

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

Plot bivariate distribution

Categorical Plots

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

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

Count Plot
`>>> sns.countplot(x="deck",
... data=titanic,`

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

Show point estimates and confidence intervals with scatterplot glyphs

Show count of observations

Point Plot
`>>> sns.pointplot(x="class",
... y="survived",
... hue="sex",
... data=titanic,
... palette="Greens_d")`

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

Show point estimates and confidence intervals as rectangular bars

Boxplot

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

Boxplot with wide-form data
Violin plot

Regression Plots

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

Plot data and a linear regression model fit

Distribution Plots

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

Plot univariate distribution

Matrix Plots

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

Heatmap

4 Further Customizations

Also see Matplotlib

Axisgrid Objects

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

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

Plot

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

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

5 Show or Save Plot

Also see Matplotlib

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

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

Close & Clear

Also see Matplotlib

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

Clear an axis
Clear an entire figure
Close a window



Python For Data Science Cheat Sheet

Matplotlib

Learn Python Interactively at www.DataCamp.com



Matplotlib

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



1 Prepare The Data

Also see [Lists & NumPy](#)

1D Data

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

2D Data or Images

```
>>> data = 2 * np.random.random((10, 10))  
>>> data2 = 3 * np.random.random((10, 10))  
>>> Y, X = np.meshgrid(np.arange(3:3:100j), -3:3:100j)  
>>> V = 1 + X * Y**2  
>>> from matplotlib.cbook import get_sample_data  
>>> img = np.load(get_sample_data('axes_grid/bivariate_normal.npy'))
```

2 Create Plot

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

Figure

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

Axes

All plotting is done with respect to an Axes. In most cases, a subplot will fit your needs. A subplot is an axes on a grid system.

```
>>> fig.add_axes()  
>>> ax1 = fig.add_subplot(221) # row-col-num  
>>> ax2 = fig.add_subplot(212)  
>>> fig3, axes = plt.subplots(nrows=2, ncols=2)  
>>> fig4, axes2 = plt.subplots(ncols=3)
```

3 Plotting Routines

1D Data

```
>>> fig, ax = plt.subplots()  
>>> lines = ax.plot(x, y)  
>>> ax.scatter(x, y)  
>>> axes[0,0].bar([1,2,3], [3,4,5])  
>>> axes[1,0].barh([0.5,1,2.5], [0,1,2])  
>>> axes[1,1].axhline(0.45)  
>>> axes[0,1].axvline(0.65)  
>>> ax.hlines(y, x_min=x1, x_max=x2)  
>>> ax.vlines(x, y_min=y1, y_max=y2, color='yellow')
```

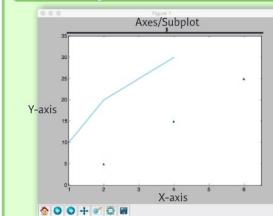
2D Data or Images

```
>>> fig, ax = plt.subplots()  
>>> im = ax.imshow(img, cmap='gist_earth',  
    interpolation='nearest',  
    vmin=-2,  
    vmax=2)
```

Colormapped or RGB arrays

Plot Anatomy & Workflow

Plot Anatomy



Workflow

The basic steps to creating plots with matplotlib are:

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

## 4 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, color='red')
>>> fg.colorbar(im, orientation='horizontal')
>>> im = ax.imshow(img,
 cmap='seismic')
```

### Markers

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

### Line Styles

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

### Text & Annotations

```
>>> ax.text(1,-2.1,
 'Example Graph',
 style='italic')
>>> ax.annotate('line',
 xy=(18, 0),
 xycoords='data',
 xytext=(10.5, 0),
 textcoords='data',
 arrowprops=dict(arrowstyle="->",
 connectionstyle="arc3"),
 rotation=15)
```

### Vector Fields

```
>>> axes[0,1].arrow(0,0,0.5,0.5)
>>> axes[1,1].quiver(y,z)
>>> axes[0,1].streamplot(X,Y,u,v)
```

### Data Distributions

```
>>> ax1.hist(y)
>>> ax3.boxplot(y)
>>> ax3.violinplot(z)
```

### MathText

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

### Limits, Legends & Layouts

```
>>> ax.margins(x=0,y=0.1)
>>> ax.axis('equal')
>>> ax.set_xlim([0,10.5],ylim=[-1.5,1.5])
>>> ax.set_xlim(0,10.5)
```

Legends

```
>>> ax.set(title='An Example Axes',
 ylabel='Y-Axis',
 xlabel='X-Axis')
>>> ax.legend(loc='best')
```

### Ticks

```
>>> ax.xaxis.set(ticks=range(1,5),
 ticklabels=[3,100,-12,"foo"])
>>> ax.tick_params(axis='y',
 direction='inout',
 length=10)
```

### Subplot Spacing

```
>>> fig3.subplots_adjust(wspace=0.5,
 hspace=0.3,
 left=0.125,
 right=0.9,
 top=0.9,
 bottom=0.1)
```

### Axis Spines

```
>>> ax1.spines['top'].set_visible(False)
>>> ax1.spines['bottom'].set_position('outward',10)
```

Add padding to a plot  
Set the aspect ratio of the plot to 1  
Set limits for x-and y-axis  
Set limits for x-axis

Set a title and x-and y-axis labels  
No overlapping plot elements

Manually set x-ticks  
Make y-ticks longer and go in and out

Adjust the spacing between subplots

Fit subplot(s) to the figure area

Make the top axis line for a plot invisible  
Move the bottom axis line outward

## 5 Save Plot

### Save Figures

```
>>> plt.savefig('foo.png')
```

### Save Transparent Figures

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

## 6 Show Plot

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

## Close & Clear

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

Clear an axis  
Clear the entire figure  
Close a window



# Python For Data Science Cheat Sheet

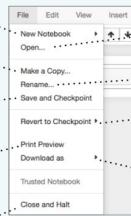
## Jupyter Notebook

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



### Saving/Loading Notebooks

Create new notebook



Make a copy of the current notebook

Open an existing notebook

Rename notebook

Save current notebook and record checkpoint

Revert notebook to a previous checkpoint

Preview of the printed notebook

Download notebook as (IPython notebook, Python, HTML, Markdown, reST, LaTeX, PDF)

Close notebook & stop running any scripts

### 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

R  
IRKernel

IJ[.]  
Julia

Installing Jupyter Notebook will automatically install the IPython kernel.

Restart kernel

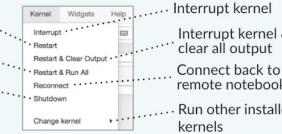
Restart kernel & run all cells

Restart kernel & run all cells

Interrupt kernel

Interrupt kernel & clear all output

Connect back to a remote notebook



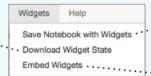
Run other installed kernels

### Widgets

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

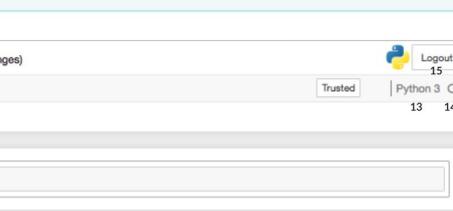
You can use them to build interactive GUIs for your notebooks or to synchronize stateful and stateless information between Python and JavaScript.

Download serialized state of all widget models in use



Save notebook with interactive widgets

Embed current widgets



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

### Writing Code And Text

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

#### Edit Cells

Cut currently selected cells to clipboard

Copy cells from clipboard to current cursor position

1. Save and checkpoint
2. Insert cell below
3. Cut cell
4. Copy cell(s)
5. Paste cell(s) below
6. Move cell up
7. Move cell down
8. Run current cell

9. Interrupt kernel

10. Restart kernel

11. Display characteristics

12. Open command palette

13. Current kernel

14. Kernel status

15. Log out from notebook server

#### Executing Cells

Walk through a UI tour

Edit the built-in keyboard shortcuts

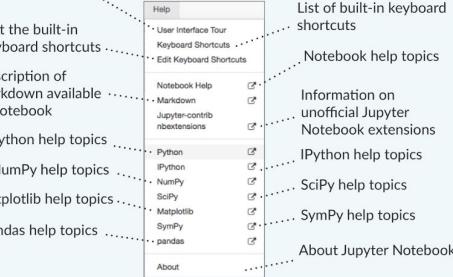
Description of markdown available in notebook

Python help topics

NumPy help topics

Matplotlib help topics

Pandas help topics



### Writing Code And Text

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

#### Edit Cells

Cut currently selected cells to clipboard

Copy cells from clipboard to current cursor position

Paste cells from clipboard above current cell

Paste cells from clipboard below current cell

Paste cells from clipboard on top of current cell

Delete current cells

Revert "Delete Cells" invocation

Split up a cell from current cursor position

Merge current cell with the one above

Merge current cell with the one below

Merge current cell with the one below

Move current cell down

Move current cell up

Find and replace in selected cells

Adjust metadata underlying the current notebook

Remove cell attachments

Paste attachments of current cell

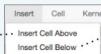
Copy attachments of current cell

Insert new cell above the current one

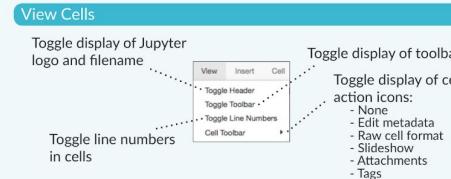
Insert image in selected cells

#### Insert Cells

Add new cell above the current one



Add new cell below the current one



# Python For Data Science Cheat Sheet

## Python Basics

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



## Variables and Data Types

### Variable Assignment

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

### Calculations With Variables

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

### Types and Type Conversion

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

## Asking For Help

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

## Strings

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

### String Operations

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

## Lists

### Also see NumPy Arrays

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

### Selecting List Elements

#### Index starts at 0

##### Subset

```
>>> my_list[1]
>>> my_list[-3]
```

##### Slice

```
>>> my_list[1:3]
>>> my_list[:1]
>>> my_list[:3]
>>> my_list[:]
```

##### Subset Lists of Lists

```
>>> my_list2[1][0]
>>> my_list2[1][:2]
```

### List Operations

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

### List Methods

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

## Libraries

### Import libraries

```
>>> import numpy
>>> import numpy as np
Selective import
>>> from math import pi
```

pandas

Data analysis

learn

Machine learning

NumPy

Scientific computing

matplotlib

2D plotting

## Install Python



ANACONDA®

Leading open data science platform  
powered by Python



spyder

Free IDE that is included  
with Anaconda



jupyter

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

## Numpy Arrays

### Also see Lists

```
>>> my_list = [1, 2, 3, 4]
```

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

### Selecting Numpy Array Elements

#### Index starts at 0

##### Subset

```
>>> my_array[1]
2
```

##### Slice

```
>>> my_array[0:2]
array([1, 2])
```

##### Subset 2D Numpy arrays

```
>>> my_2darray[:,0]
array([1, 4])
```

my\_2darray[rows, columns]

### Numpy Array Operations

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

### Numpy Array Functions

```
>>> my_array.shape
>>> np.append(other_array)
>>> np.insert(my_array, 1, 5)
>>> np.delete(my_array, [1])
>>> np.mean(my_array)
>>> np.median(my_array)
>>> my_array.corrcoef()
>>> np.std(my_array)
```

Get the dimensions of the array

Append items to an array

Insert items in an array

Delete items in an array

Mean of the array

Median of the array

Correlation coefficient

Standard deviation

### String Operations

#### Index starts at 0

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

### String Methods

```
>>> my_string.upper()
>>> my_string.lower()
>>> my_string.count('w')
>>> my_string.replace('e', 'i')
>>> my_string.strip()
```

String to uppercase  
String to lowercase  
Count String elements  
Replace String elements  
Strip whitespaces

DataCamp

Learn Python for Data Science Interactively



# Python For Data Science Cheat Sheet

## PySpark - SQL Basics

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



### 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())))
>>> schemaString = "name age"
>>> fields = [StructTypeField(field.name, StringType(), True) for field.name in schemaString.split()]
>>> schema = StructType(fields)
>>> spark.createDataFrame(people, schema).show()
+-----+
| name|age|
+---+---+
| Mine| 28|
| Flin| 39|
| Jordan| 30|
+-----+
```

#### From Spark Data Sources

```
JSON
>>> df = spark.read.json("customer.json")
>>> df.show()
+-----+-----+-----+-----+
| address|firstName|lastName| phoneNumber|
+-----+-----+-----+-----+
| [New York, 10021, N...| 25| John| Smith|[212 555-1234, ho...|
| [New York, 10021, N...| 21| Jane| Doe|[322 888-1234, ho...|
+-----+-----+-----+-----+
>>> df2 = spark.read.load("people.json", format="json")
Parquet files
>>> df3 = spark.read.load("users.parquet")
TXT files
>>> df4 = spark.read.text("people.txt")
```

### 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 the first n rows               |
| >>> df.take(2) | Return the schema of df               |

### Duplicate Values

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

### Queries

```
>>> from pyspark.sql import functions as F
Select
>>> df.select("firstName").show()
>>> df.select("firstName", "lastName") \
 .show()
>>> df.select("firstName",
 "age",
 explode("phoneNumber") \
 .alias("contactInfo")) \
 .select("contactInfo.type",
 "firstName",
 "age") \
 .show()
>>> df.select(df["firstName"], df["age"] + 1) \
 .show()
>>> df.select(df["age"] > 24).show()
When
>>> df.select("firstName",
 F.when(df.age > 30, 1) \
 .otherwise(0)) \
 .show()
>>> df[df.firstName.isin("Jane", "Boris")] \
 .collect()
Like
>>> df.select("firstName",
 df.lastName.like("Smith")) \
 .show()
Startswith - Endswith
>>> df.select("firstName",
 df.lastName \
 .startswith("Sm")) \
 .show()
>>> df.select(df.lastName.endswith("th")) \
 .show()
Substring
>>> df.select(df.firstName.substr(1, 3)) \
 .alias("name") \
 .collect()
Between
>>> df.select(df.age.between(22, 24)) \
 .show()
```

Show all entries in firstName column

Show all entries in firstName, age and type

Show all entries in firstName and age, add 1 to the entries of age

Show all entries where age >24

Show firstName and 0 or 1 depending on age >30

Show firstName if in the given options

Show firstName, and lastName is TRUE if lastName is like Smith

Show firstName, and TRUE if lastName starts with Sm

Show last names ending in th

Return substrings of firstName

Show age: values are TRUE if between 22 and 24

### 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("telephoneType",
 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)
```

### GroupBy

```
>>> df.groupBy("age") \
 .count() \
 .show()
```

Group by age, count the members in the groups

### Filter

```
>>> df.filter(df["age"]>24).show()
```

Filter entries of age, only keep those records of which the values are >24

### Sort

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

### Missing & Replacing Values

```
>>> df.na.fill(50).show()
>>> df.na.drop().show()
>>> df.na \
 .replace(10, 20) \
 .show()
```

Replace null values  
Return new df omitting rows with null values  
Return new df replacing one value with another

### Repartitioning

```
>>> df.repartition(10) \
 .rdd \
 .getNumPartitions()
>>> df.coalesce(1).rdd.getNumPartitions()
```

df with 10 partitions  
df with 1 partition

### Running SQL 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()
```

### Output

#### Data Structures

```
>>> rdd = df.rdd
>>> df.toJSON().first()
>>> df.toPandas()
```

Convert df into an RDD  
Convert df into a RDD of string  
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()
```

