

# Data Wrangling

with pandas Cheat Sheet  
<http://pandas.pydata.org>

[Pandas API Reference](#) [Pandas User Guide](#)

## 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.

		a	b	c
N	v			
D	1	4	7	10
	2	5	8	11
e	2	6	9	12

```
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

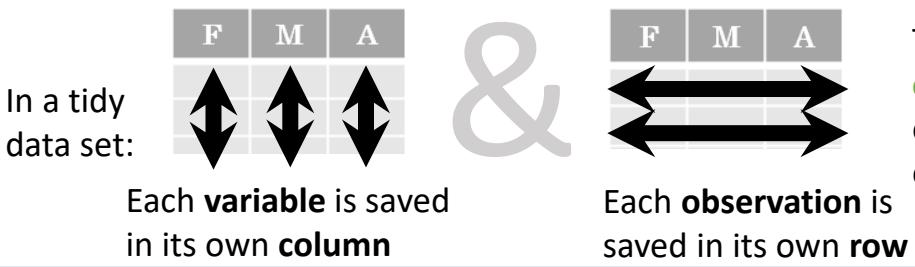
## Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result.

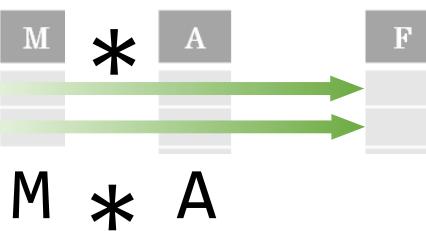
This improves readability of code.

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

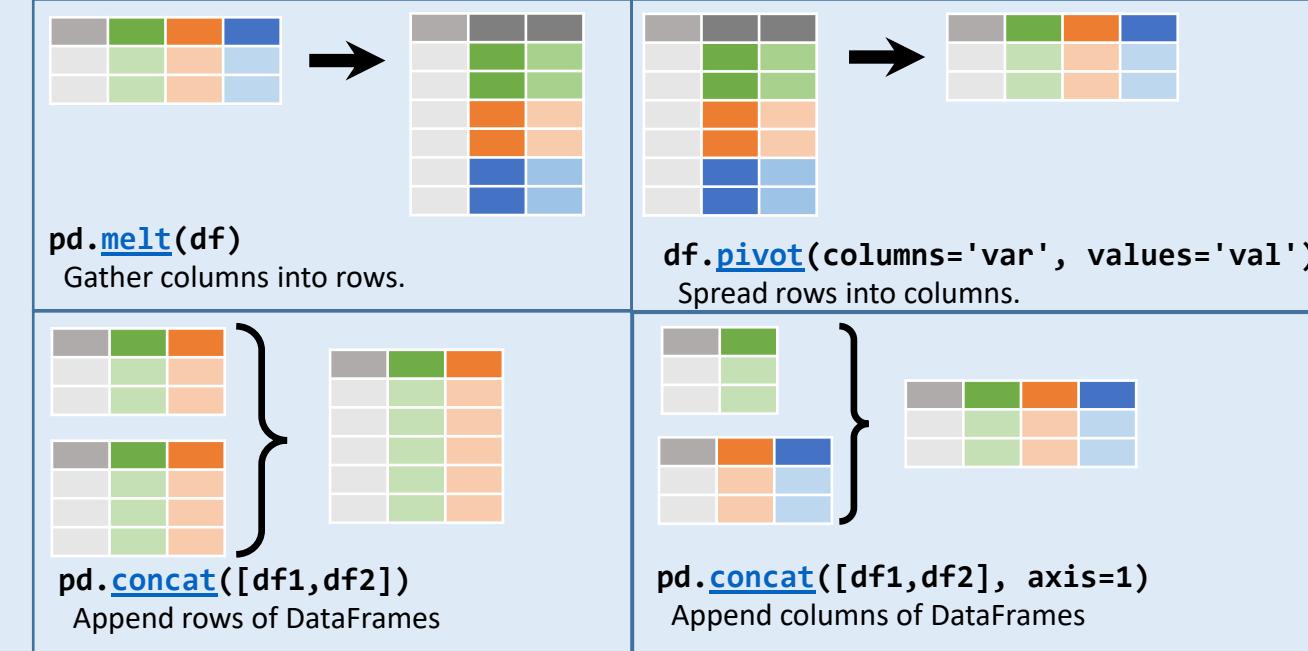
# Tidy Data – A foundation for wrangling in pandas



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



## Reshaping Data – Change layout, sorting, reindexing, renaming



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

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

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

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

`df.head(n)`  
Select first n rows.

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

## 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.

## Using query

`query()` allows Boolean expressions for filtering rows.

`df.query('Length > 7')`  
`df.query('Length > 7 and Width < 8')`  
`df.query('Name.str.startswith("abc")', engine="python")`

Use `df.loc[]` and `df.iloc[]` to select only rows, only columns or both.

Use `df.at[]` and `df.iat[]` to access a single value by row and column.  
First index selects rows, second index columns.

`df.iloc[10:20]`  
Select rows 10-20.

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

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

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

`df.iat[1, 2]` Access single value by index  
`df.at[4, 'A']` Access single value by label

## Logic in Python (and pandas)

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

## 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'

# Summarize Data

`df['w'].value_counts()`

Count number of rows with each unique value of variable

`len(df)`

# of rows in DataFrame.

`df.shape`

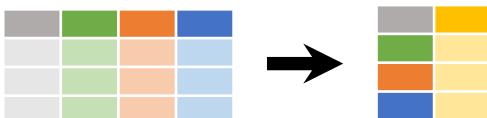
Tuple of # of rows, # of columns in DataFrame.

`df['w'].nunique()`

# of distinct values in a column.

`df.describe()`

Basic descriptive and 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

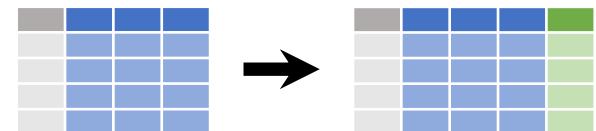
`df.dropna()`

Drop rows with any column having NA/null data.

`df.fillna(value)`

Replace all NA/null data with value.

# Make New Columns



`df.assign(Area=lambda df: df.Length*df.Height)`

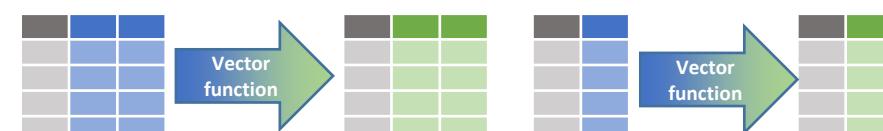
Compute and append one or more new columns.

`df['Volume'] = df.Length*df.Height*df.Depth`

Add single column.

`pd.qcut(df.col, n, labels=False)`

Bin column into n buckets.



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

`max(axis=1)`

Element-wise max.

`clip(lower=-10,upper=10)`

Trim values at input thresholds

`min(axis=1)`

Element-wise min.

`abs()`

Absolute value.

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

`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.

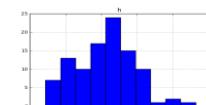
`cumprod()`

Cumulative product.

# Plotting

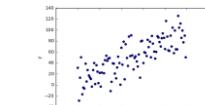
`df.plot.hist()`

Histogram for each column



`df.plot.scatter(x='w',y='h')`

Scatter chart using pairs of points



# Combine Data Sets

`adf`

x1	x2
A	1
B	2
C	3

`bdf`

x1	x3
A	T
B	F
D	T



Standard Joins

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

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

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

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

x1	x2	x3
A	1	T
B	2	F

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

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

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

Filtering Joins

x1	x2
A	1
B	2

x1	x2
C	3

`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`

x1	x2
A	1
B	2
C	3

`zdf`

x1	x2
B	2

# Matplotlib for beginners

Matplotlib is a library for making 2D plots in Python. It is designed with the philosophy that you should be able to create simple plots with just a few commands:

## 1 Initialize

```
import numpy as np  
import matplotlib.pyplot as plt
```

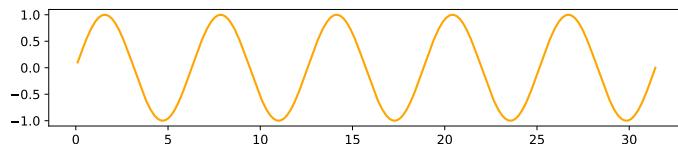
## 2 Prepare

```
X = np.linspace(0, 4*np.pi, 1000)  
Y = np.sin(X)
```

## 3 Render

```
fig, ax = plt.subplots()  
ax.plot(X, Y)  
plt.show()
```

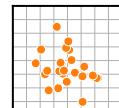
## 4 Observe



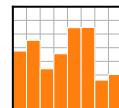
## Choose

Matplotlib offers several kind of plots (see Gallery):

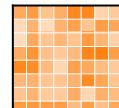
```
X = np.random.uniform(0, 1, 100)  
Y = np.random.uniform(0, 1, 100)  
ax.scatter(X, Y)
```



```
X = np.arange(10)  
Y = np.random.uniform(1, 10, 10)  
ax.bar(X, Y)
```



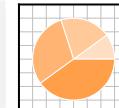
```
Z = np.random.uniform(0, 1, (8,8))  
ax.imshow(Z)
```



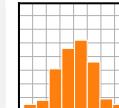
```
Z = np.random.uniform(0, 1, (8,8))  
ax.contourf(Z)
```



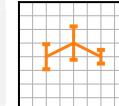
```
Z = np.random.uniform(0, 1, 4)  
ax.pie(Z)
```



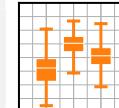
```
Z = np.random.normal(0, 1, 100)  
ax.hist(Z)
```



```
X = np.arange(5)  
Y = np.random.uniform(0, 1, 5)  
ax.errorbar(X, Y, Y/4)
```



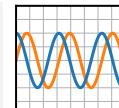
```
Z = np.random.normal(0, 1, (100,3))  
ax.boxplot(Z)
```



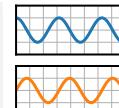
## Organize

You can plot several data on the same figure, but you can also split a figure in several subplots (named Axes):

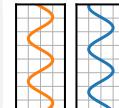
```
X = np.linspace(0, 10, 100)  
Y1, Y2 = np.sin(X), np.cos(X)  
ax.plot(X, Y1, X, Y2)
```



```
fig, (ax1, ax2) = plt.subplots(2,1)  
ax1.plot(X, Y1, color="C1")  
ax2.plot(X, Y2, color="C0")
```



```
fig, (ax1, ax2) = plt.subplots(1,2)  
ax1.plot(Y1, X, color="C1")  
ax2.plot(Y2, X, color="C0")
```

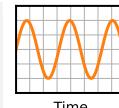


## Label (everything)

```
ax.plot(X, Y)  
fig.suptitle(None)  
ax.set_title("A Sine wave")
```



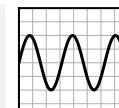
```
ax.plot(X, Y)  
ax.set_ylabel(None)  
ax.set_xlabel("Time")
```



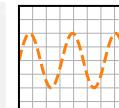
## Tweak

You can modify pretty much anything in a plot, including limits, colors, markers, line width and styles, ticks and ticks labels, titles, etc.

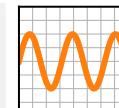
```
X = np.linspace(0, 10, 100)  
Y = np.sin(X)  
ax.plot(X, Y, color="black")
```



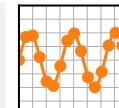
```
X = np.linspace(0, 10, 100)  
Y = np.sin(X)  
ax.plot(X, Y, linestyle="--")
```



```
X = np.linspace(0, 10, 100)  
Y = np.sin(X)  
ax.plot(X, Y, linewidth=5)
```



```
X = np.linspace(0, 10, 100)  
Y = np.sin(X)  
ax.plot(X, Y, marker="o")
```



## Explore

Figures are shown with a graphical user interface that allows to zoom and pan the figure, to navigate between the different views and to show the value under the mouse.

## Save (bitmap or vector format)

```
fig.savefig("my-first-figure.png", dpi=300)  
fig.savefig("my-first-figure.pdf")
```

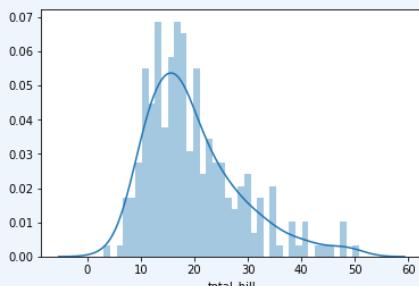
# Cheatography

## Python - Seaborn IMGs Cheat Sheet by DarioPittera (aggialavura) via [cheatography.com/83764/cs/19859/](https://cheatography.com/83764/cs/19859/)

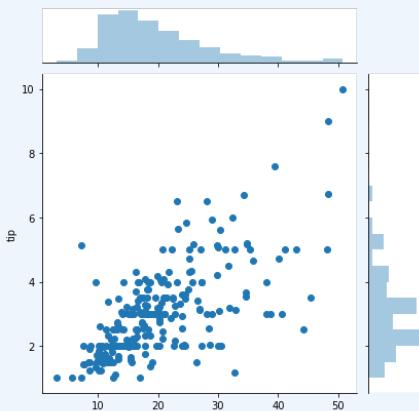
### TO START

```
import seaborn as sns  
%matplotlib inline
```

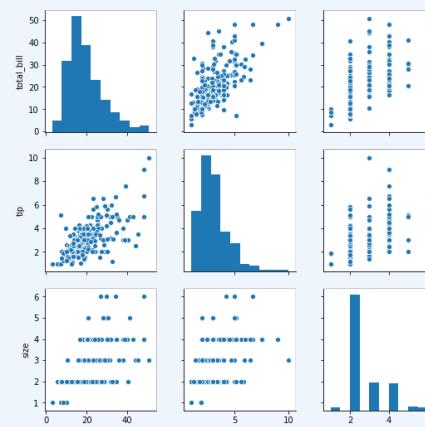
### distplot()



### jointplot()

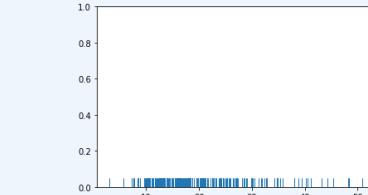


### pairplot()

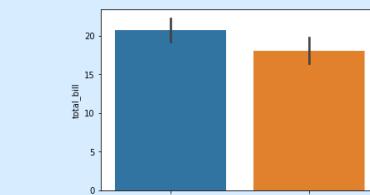


```
sns.pairplot(tips)
```

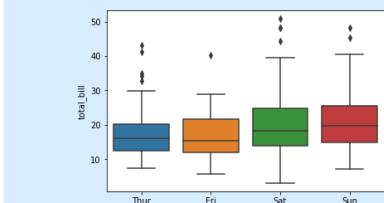
### rugplot()



### barplot()

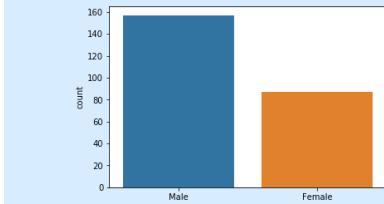


### boxplot()



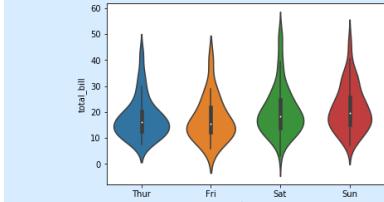
```
sns.boxplot(x='day',y='total_bill',data = tips)
```

### countplot()



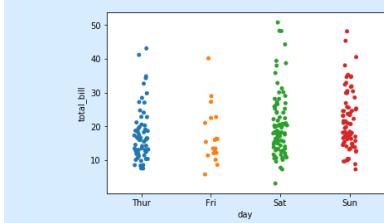
```
sns.countplot(x='sex',data=tips)
```

### violinplot()



```
sns.violinplot(x="day", y="total_bill", data=tips)
```

### stripplot()



```
sns.stripplot(x="day", y="total_bill", data=tips)
```

Sponsored by [Readable.com](https://readable.com)

Measure your website readability!

<https://readable.com>



By DarioPittera (aggialavura)

[cheatography.com/aggialavura/](https://cheatography.com/aggialavura/)  
[www.dariopittera.com](http://www.dariopittera.com)

Not published yet.

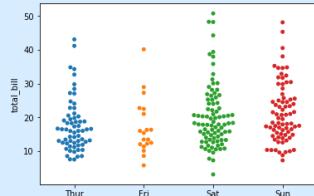
Last updated 24th June, 2019.

Page 1 of 3.

# Cheatography

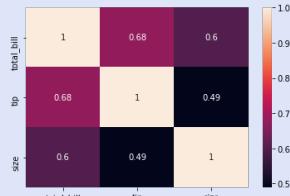
Python - Seaborn IMGs Cheat Sheet  
by DarioPittera (aggialavura) via [cheatography.com/83764/cs/19859/](http://cheatography.com/83764/cs/19859/)

## swarmplot()



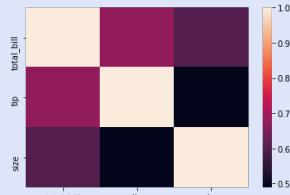
```
sns.swarmplot(x="day", y="total_bill",  
               data=tips)
```

## heatmap() with values



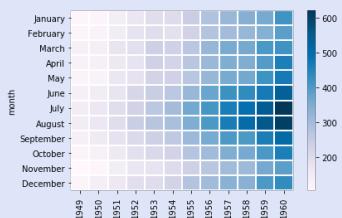
```
sns.heatmap(tips.corr(), annot=True)
```

## heatmap()



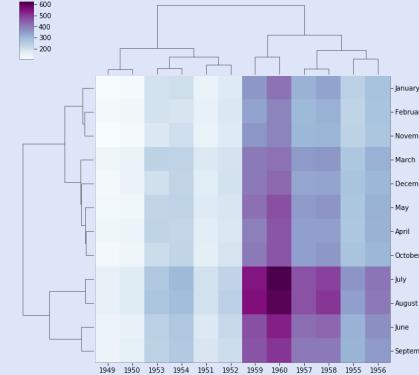
```
sns.heatmap(tips.corr())
```

## heatmap(linecolor=' ', linewidth=x')



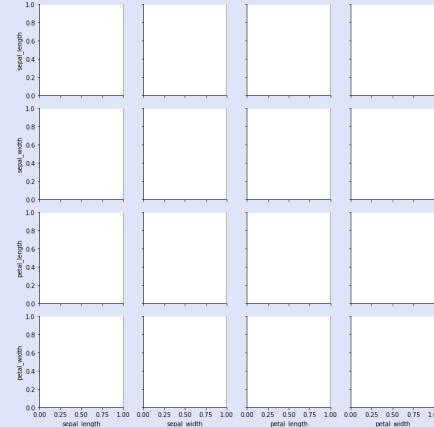
```
sns.heatmap(matrix, cmap='PuBu', linecolor='white', linewidths=1)
```

## clustermap()



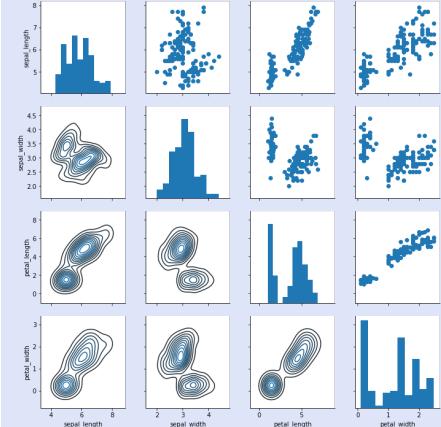
```
sns.clustermap(matrix, cmap='BuPu')
```

## PairGrid()



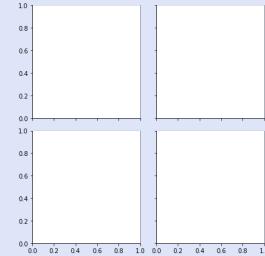
```
sns.PairGrid(df)
```

## PairGrid() with options



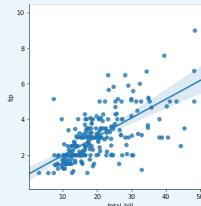
```
g = sns.PairGrid(df)  
g.map_diag(plt.hist)  
g.map_upper(plt.scatter)  
g.map_lower(sns.kdeplot)
```

## FacetGrid()



```
g = sns.FacetGrid(data=tips, col='time',  
                   row='smoker')
```

## lmplot()



```
sns.lmplot(x='total_bill', y='tip', data=tips)
```

Not published yet.

Last updated 24th June, 2019.

Page 2 of 3.

Sponsored by [Readable.com](https://readable.com)

Measure your website readability!

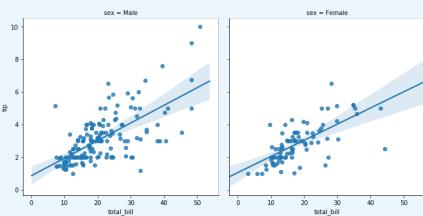
<https://readable.com>

# Cheatography

## Python - Seaborn IMGs Cheat Sheet

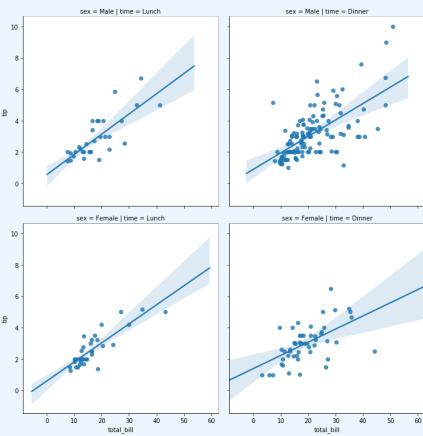
by DarioPittera (aggialavura) via [cheatography.com/83764/cs/19859/](https://cheatography.com/83764/cs/19859/)

### lmplot(col="")



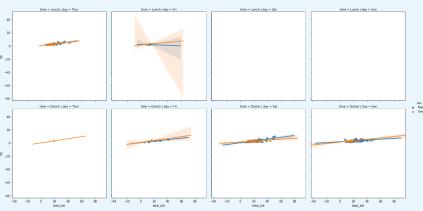
```
sns.lmplot(x="total_bill",y="tip",data=tips,c-  
ol="sex")
```

### lmplot(col="", row="")



```
sns.lmplot(x="total_bill", y="tip", row="sex",  
col="time",data=tips)
```

### lmplot(col="", row="", hue="")



```
sns.lmplot(x="total_bill",y="tip",data=tips,c-  
ol='day', row='time', hue='sex')
```



By DarioPittera (aggialavura)

[cheatography.com/aggialavura/](https://cheatography.com/aggialavura/)  
[www.dariopittera.com](http://www.dariopittera.com)

Not published yet.

Last updated 24th June, 2019.

Page 3 of 3.

Sponsored by **Readable.com**

Measure your website readability!

<https://readable.com>

# Python For Data Science

## Scikit-Learn Cheat Sheet

Learn Scikit-Learn online at [www.DataCamp.com](http://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[:, :2], iris.target
>>> X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=33)
>>> scaler = preprocessing.StandardScaler().fit(X_train)
>>> X_train = scaler.transform(X_train)
>>> X_test = scaler.transform(X_test)
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
>>> knn.fit(X_train, y_train)
>>> y_pred = knn.predict(X_test)
>>> accuracy_score(y_test, y_pred)
```

### > Loading The Data

Also see NumPy & Pandas

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

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

### > Training And Test Data

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

### > Model Fitting

**Supervised learning**

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

**Unsupervised Learning**

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

### > Prediction

**Supervised Estimators**

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

**Unsupervised Estimators**

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

### > Preprocessing The Data

#### Standardization

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

#### Normalization

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

#### Binarization

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

#### Encoding Categorical Features

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

#### Imputing Missing Values

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

#### Generating Polynomial Features

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

### > Create Your Model

#### Supervised Learning Estimators

##### Linear Regression

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

##### Support Vector Machines (SVM)

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

##### Naive Bayes

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

##### KNN

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

#### Unsupervised Learning Estimators

##### Principal Component Analysis (PCA)

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

##### K Means

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

### > Evaluate Your Model's Performance

#### Classification Metrics

##### Accuracy Score

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

##### Classification Report

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

##### Confusion Matrix

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

#### Regression Metrics

##### Mean Absolute Error

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

##### Mean Squared Error

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

##### R<sup>2</sup> Score

```
>>> from sklearn.metrics import r2_score
>>> r2_score(y_true, y_pred)
```

#### Clustering Metrics

##### Adjusted Rand Index

```
>>> from sklearn.metrics import adjusted_rand_score
>>> adjusted_rand_score(y_true, y_pred)
```

##### Homogeneity

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

##### V-measure

```
>>> from sklearn.metrics import v_measure_score
>>> metrics.v_measure_score(y_true, y_pred)
```

#### Cross-Validation

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

### > Tune Your Model

#### Grid Search

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

#### Randomized Parameter Optimization

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

### Basic graph manipulation

```
import networkx as nx  
G=nx.Graph()
```

```
G=nx.MultiGraph()
```

```
G.add_edges_from([(0, 1),(0, 2),(1, 3),(2, 4)])
```

```
nx.drawnetworkx(G)
```

```
G.add_node('A', role='manager')
```

```
G.add_edge('A', 'B', relation = 'friend')
```

```
G.nodes()['A']['role'] = 'team member'
```

```
G.nodes()['A'], G.edges()['A', 'B'])
```

```
G.edges(), G.nodes()
```

```
list(G.edges())
```

```
G.nodes(data=True), G.edges(data=True)
```

```
G.nodes(data='relation')
```

### Creating graphs from data

### Bipartite graphs (cont)

```
bipartite.sets(B)
```

Create a graph allowing parallel edges

Create graph from edges

Draw the graph

Add a node

Add an edge

Set attribute of a node

View attributes of node, edge

Show edges, nodes

Return as list instead of EdgeView class

Include node/edge attributes

Return specific attribute

### Network Connectivity

Get each set of nodes of bipartite graph

Bipartite projected graph - nodes with bipartite friends in common

Projected graph with weights (number of friends in common)

G=nx.read_edgelist('G_edgelist.txt', nodetype=int)	<code>Create from nx.adjacency_matrix(G) list</code>	Local clustering coefficient
G=nx.Graph(G_matrix)	<code>Create from nx.transitivity(G)</code>	Global clustering coefficient
G=nx.read_edgelist('G_edgelist.txt', data=[('Weight', nx.shortest_path_length(G, n1, n2))], int))	<code>Create from nx.shortest_path(G, n1, n2)</code>	Transitivity (% of open triads)
G=nx.from_pandas_dataframe(pd.read_csv('G_df', 'n1', 'n2', edge_attr='weight'))	<code>Create from df</code>	Outputs the path itself
T=nx.bfs_tree(G, n1)	<code>Create from nx.average_shortest_path_length(G)</code>	Create breadth-first search tree from node n1
Adjacency list format: 0 1 2 3 5 1 3 6 ...	<code>nx.diameter(G)</code>	Average distance between all pairs of nodes
Edgelist format: 0 1 14 0 2 17	<code>nx.eccentricity(G)</code>	Maximum distance between any pair of nodes

## Bipartite graphs

from networkx.algorithms import bipartite		
bipartite.is_bipartite(B)	Check if graph B is bipartite	
bipartite.is_bipartition_nodes_set(B, set)	Check if set of nodes is bipartition of graph	

<code>nx.eccentricity(G)</code>	Returns each node's distance to furthest node
<code>nx.radius(G)</code>	Minimum eccentricity in the graph
<code>nx.periphery(G)</code>	Set of nodes where eccentricity=diameter
<code>nx.centrality(G)</code>	Set of nodes where eccentricity=radius

## Connectivity: Network Robustness

<code>nx.nodes_connectivity(G)</code>	Min nodes removed to disconnect a network
<code>nx.minimum_node_cut()</code>	Which nodes?
<code>nx.edges_connectivity(G)</code>	Min edges removed to disconnect a network
<code>nx.minimum_edge_cut(G)</code>	Which edges?
<code>nx.all_simple_paths(G, n1, n2)</code>	Show all paths between two nodes



By **RJ Murray** (murenei)  
[cheatography.com/murenei/](http://cheatography.com/murenei/)  
[tutify.com.au](http://tutify.com.au)

Published 4th June, 2018.  
 Last updated 4th June, 2018.  
 Page 1 of 5.

Sponsored by **Readable.com**  
 Measure your website readability!  
<https://readable.com>

### Network Connectivity: Connected Components

<code>nx.is_connected(G)</code>	Is there a path between every pair of nodes?
<code>nx.number_of_connected_components(G)</code>	# separate components
<code>nx.node_connected_component(N)</code>	Which connected component does $N$ belong to?
<code>nx.is_strongly_connected(G)</code>	Is the network connected directionally?
<code>nx.is_weakly_connected(G)</code>	Is the directed network connected if assumed undirected?

### Influence Measures and Network Centralization (cont)

<code>nx.edge_betweenness_centrality(G)</code>	B e e s o e
<code>nx.edge_betweenness_centrality_subset(G, s)</code>	s o e

Normalization: Divide by number of pairs of nodes.

### PageRank and Hubs & Authorities Algorithms

<code>nx.pagerank(G, alpha=0.8)</code>	Scaled PageRank of $G$ with dampening parameter
<code>h, a=nx.hits(G)</code>	HITS algorithm - outputs 2 dictionaries (hubs, authorities)
<code>h, a=nx.hits(G, max_iter=10, normализed=True)</code>	Constrained HITS and normalized by sum at each stage

Centrality measures make different assumptions about what it means to be a “central” node. Thus, they produce different rankings.

### Network Evolution - Real-world Applications

<code>G.degree(), G.in_degree(), G.out_degree()</code>	Distribution of node degrees
Preferential Attachment Model	Results in power law -> many nodes with low degrees; few with high degrees

### Common Graphs

<code>G=nx.karate_club_graph()</code>	Karate club graph (social network)
<code>G=nx.path_graph(n)</code>	Path graph with $n$ nodes
<code>G=nx.complete_graph(n)</code>	Complete graph on $n$ nodes
<code>G=random_regular_graph(d, n)</code>	Random $d$ -regular graph on $n$ -nodes

See [NetworkX Graph Generators reference](#) for more.

Also see “An Atlas of Graphs” by Read and Wilson (1998).

### Influence Measures and Network Centralization

dc=nx.degree_centrality(G)	Degree centrality for network	Preferential Attachment Model with $n$ nodes and each new node attaching to $m$ existing nodes
dc[node]	Degree centrality for a node	
nx.in_degree_centrality(G), nx.out_degree_centrality(G)	DC for directed networks	
cc=nx.closeness_centrality(G, normalized=True)	Closeness centrality Small World model (normalised) for the network	High average degree (global clustering) and low average shortest path
cc[node]	Closeness centrality for an individual	
bC=nx.betweenness_centrality(..., normalized=True, ..., endpoints=False, ...)	Normalized betweenness centrality BC excluding endpoints	Small World network of $n$ nodes, connected to its $k$ nearest neighbours, with chance $p_c$ rewiring
..., K=10, ...)	BC approximated using random sample of $K$ nodes	$t = \max$ iterations to try to ensure connected graph
nx.betweenness_centrality_subset(G, {subset})	BC calculated on subset	$p = \text{probability of adding (no rewiring)}$
Link Prediction measures		How likely are 2 nodes to connect, given an existing network



By **RJ Murray** (murenei)  
[cheatography.com/murenei/](http://cheatography.com/murenei/)  
[tutify.com.au](http://tutify.com.au)

Published 4th June, 2018.  
Last updated 4th June, 2018.  
Page 2 of 5.

Sponsored by **Readable.com**  
Measure your website readability!  
<https://readable.com>

### Network Evolution - Real-world Applications (cont)

<code>nx.common_neighor_s(G, n1, n2)</code>	Calc common neighbors of nodes $n_1, n_2$
<code>nx.jaccard_coefficient(G)</code>	Normalised common neighbors measure
<code>nx.resource_allocation_index(G)</code>	Calc RAI of all nodes not already connected by an edge
<code>nx.adamic_adar_index(G)</code>	As per RAI but with log of degree of common neighbor
<code>nx.product_degree(G)</code>	Product of two nodes' degrees
Community Common Neighbors	Common neighbors but with bonus if they belong in same 'community'
<code>nx.community_common_neighbors(n1, n2)</code>	CCN score for $n_1, n_2$
<code>G.node['A']['community']=1</code>	Add community attribute to node
<code>nx.resource_allocation_index(G)</code>	Community Resource Allocation score

These scores give only an indication of whether 2 nodes are likely to connect.

To make a link prediction, you would use these scores as features in a classification ML model.



By **RJ Murray** (murenei)  
[cheatography.com/murenei/](https://cheatography.com/murenei/)  
[tutify.com.au](https://tutify.com.au)

Published 4th June, 2018.  
Last updated 4th June, 2018.  
Page 3 of 5.

Sponsored by **Readable.com**  
Measure your website readability!  
<https://readable.com>

# Python For Data Science Cheat Sheet

## NumPy Basics

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



### NumPy

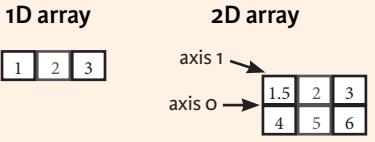
The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.

Use the following import convention:

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



### NumPy Arrays



### Creating Arrays

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

### Initial Placeholders

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

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

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

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

### I/O

#### Saving & Loading On Disk

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

#### Saving & Loading Text Files

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

### Data Types

```
>>> np.int64
>>> np.float32
>>> np.complex
>>> np.bool
>>> np.object
>>> np.string_
>>> np_unicode_
```

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

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

### Asking For Help

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

### Array Mathematics

#### Arithmetic Operations

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

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

### Comparison

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

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

### Aggregate Functions

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

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

### Copying Arrays

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

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

### Sorting Arrays

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

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

### Subsetting, Slicing, Indexing

#### Subsetting

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

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

#### Slicing

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

Select items at index 0 and 1  
Select items at rows 0 and 1 in column 1

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

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

```
>>> a[ : :-1]
      array([3, 2, 1])
```

Reversed array `a`

#### Boolean Indexing

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

Select elements from `a` less than 2

#### Fancy Indexing

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

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

Also see Lists

Select the element at the 2nd index

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

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[0:1, :]`)  
Same as `[1, :, :]`

Reversed array `a`

Select elements from `a` less than 2

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

Flatten the array  
Reshape, but don't change data

#### Adding/Removing Elements

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

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

#### Combining Arrays

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

Concatenate arrays

Stack arrays vertically (row-wise)  
Stack arrays vertically (row-wise)

Stack arrays horizontally (column-wise)  
Create stacked column-wise arrays

Create stacked column-wise arrays

#### Splitting Arrays

```
>>> np.hsplit(x,3)
      [array([1]),array([2]),array([3])]
>>> np.vsplit(c,2)
      [array([[ 1.5,  2.,  3.],
              [ 4.,  5.,  6.]]),
       array([[ 3.,  2.,  1.],
              [ 4.,  5.,  6.]])]
```

Split the array horizontally at the 3rd index  
Split the array vertically at the 2nd index



Data Types	
Integer	-256, 15
Float	-253.23, 1.253e-10
String	"Hello", 'Goodbye', """Multiline""""
Boolean	True, False
List	[ value, ... ]
Tuple	( value, ... ) <sup>1</sup>
Dictionary	{ key: value, ... }
Set	{ value, value, ... } <sup>2</sup>
<sup>1</sup> Parentheses usually optional	
<sup>2</sup> Create an empty set with set()	

Statements	
If Statement	
if expression:	statements
elif expression:	statements
else:	statements
While Loop	
while expression:	statements
For Loop	
for var in collection:	statements
Counting For Loop	
for i in range(start, end [, step]):	statements (start is included; end is not)

Arithmetic Operators			
x + y	add	x - y	subtract
x * y	multiply	x / y	divide
x % y	modulus	x ** y	x <sup>y</sup>

Assignment shortcuts: x op= y  
Example: x += 1 increments x

Comparison Operators			
x < y	Less	x <= y	Less or eq
x > y	Greater	x >= y	Greater or eq
x == y	Equal	x != y	Not equal

Boolean Operators		
not x	x and y	x or y

Exception Handling	
try:	
statements	
except [ exception type [ as var ] ]:	statements
finally:	statements

Conversion Functions	
int(expr)	Converts expr to integer
float(expr)	Converts expr to float
str(expr)	Converts expr to string
chr(num)	ASCII char num

String / List / Tuple Operations	
len(s)	length of s
s[i]	i <sup>th</sup> item in s (0-based)
s[start : end]	slice of s from start (included) to end (excluded)
x in s	True if x is contained in s
x not in s	True if x is not contained in s
s + t	the concatenation of s with t
s * n	n copies of s concatenated
sorted(s)	a sorted copy of s
)	
s.index( item)	position in s of item

More String Operations	
s.lower()	lowercase copy of s
s.replace(old, new)	copy of s with old replaced with new
s.split( delim )	list of substrings delimited by delim

More String Operations (cont)	
s.strip()	copy of s with whitespace trimmed
s.upper()	uppercase copy of s
See also	
<a href="http://docs.python.org/library/stdtypes.html#string-methods">http://docs.python.org/library/stdtypes.html#string-methods</a>	

Mutating List Operations	
del lst[i]	Deletes i <sup>th</sup> item from lst
lst.append(e)	Appends e to lst
lst.insert(i, e)	Inserts e before i <sup>th</sup> item in lst
lst.sort()	Sorts lst
See also	
<a href="http://docs.python.org/library/stdtypes.html#typeseq-mutable">http://docs.python.org/library/stdtypes.html#typeseq-mutable</a>	

Dictionary Operations	
len(d)	Number of items in d
del d[key]	Removes key from d
key in d	True if d contains key
d.keys()	Returns a list of keys in d
See also	
<a href="http://docs.python.org/library/stdtypes.html#mapping-types-dict">http://docs.python.org/library/stdtypes.html#mapping-types-dict</a>	

Function Definitions	
def name(arg1, arg2, ...):	statements
return expr	
Environment	
sys.argv	List of command line arguments (argv[0] is executable)
os.environ	Dictionary of environment variables
os.curdir	String with path of current directory
import sys; print(sys.argv)	or
from sys import argv; print(argv)	



By sschaub  
[cheatography.com/sschaub/](http://cheatography.com/sschaub/)

Published 21st May, 2012.  
Last updated 2nd June, 2014.  
Page 1 of 2.

Sponsored by **CrosswordCheats.com**  
Learn to solve cryptic crosswords!  
<http://crosswordcheats.com>

### String Formatting

```
"Hello, {0} {1}".format("abe", "jones")
Hello, abe jones
"Hello, {fn} {ln}".format(fn="abe", ln="jones")
Hello, abe jones
"You owe me ${0:,.2f}".format(253422.3)
You owe me $253,422.30
now = datetime.now()
'{%Y-%m-%d %H:%M:%S}'.format(now)
2012-05-16 15:04:33
```

See also <http://docs.python.org/library/string.html#format-specification-mini-language>

### Useful Functions

<code>exit( code )</code>	Terminate program with exit <i>code</i>
<code>raw_input("prompt")</code>	Print <i>prompt</i> and readline() from stdin <sup>1</sup>

<sup>1</sup> Use `input("prompt")` in Python 3

### Code Snippets

#### Loop Over Sequence

```
for index, value in enumerate(seq):
    print("{} : {}".format(index, value))
```

#### Loop Over Dictionary

```
for key in sorted(dict):
    print(dict[key])
```

#### Read a File

```
with open("filename", "r") as f:
    for line in f:
        line = line.rstrip("\n") # Strip newline
        print(line)
```

### Other References

<http://rgruet.free.fr/>  
*Great Python 2.x Quick Reference*  
<http://www.cheatography.com/davechild/cheat-sheets/python/>  
*More Python Cheatsheet Goodness*



By **sschaub**  
[cheatography.com/sschaub/](http://cheatography.com/sschaub/)

Published 21st May, 2012.  
Last updated 2nd June, 2014.  
Page 2 of 2.

Sponsored by **CrosswordCheats.com**  
Learn to solve cryptic crosswords!  
<http://crosswordcheats.com>