

## Quick start

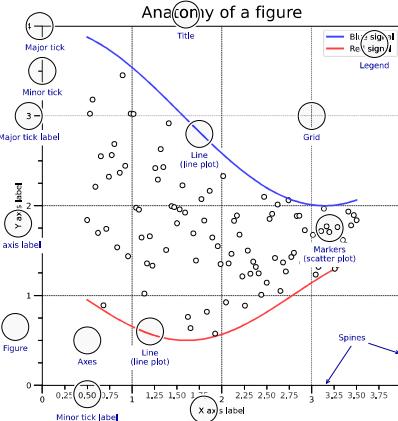
```
import numpy as np
import matplotlib as mpl
import matplotlib.pyplot as plt

X = np.linspace(0, 2*np.pi, 100)
Y = np.cos(X)

fig, ax = plt.subplots()
ax.plot(X, Y, color='green')

fig.savefig("figure.pdf")
fig.show()
```

## Anatomy of a figure



## Subplots layout

```
subplot[s](rows,cols,...) API
fig, axs = plt.subplots(3, 3)

G = gridspec(rows,cols,...) API
ax = G[0,:]

ax.inset_axes(extent) API

d=make_axes_locatable(ax) API
ax = d.new_horizontal('10%')
```

## Getting help

- [matplotlib.org](http://matplotlib.org)
- [github.com/matplotlib/matplotlib/issues](https://github.com/matplotlib/matplotlib/issues)
- [discourse.matplotlib.org](https://discourse.matplotlib.org)
- [stack overflow.com/questions/tagged/matplotlib](https://stackoverflow.com/questions/tagged/matplotlib)
- [gitter.im/matplotlib](https://gitter.im/matplotlib)
- [twitter.com/matplotlib](https://twitter.com/matplotlib)
- [Matplotlib users mailing list](mailto:Matplotlib users mailing list)

## Basic plots

```
plot([X],Y,[fmt],...) API
X, Y, fmt, color, marker, linestyle

scatter(X,Y,...) API
X, Y, [s]izes, [c]olors, marker, cmap
```

```
bar[h](x,height,...) API
x, height, width, bottom, align, color
```

```
imshow(Z,...) API
Z, cmap, interpolation, extent, origin
```

```
contour(f)([X],[Y],Z,...) API
X, Y, Z, levels, colors, extent, origin
```

```
pcolormesh([X],[Y],Z,...) API
X, Y, Z, vmin, vmax, cmap
```

```
quiver([X],[Y],U,V,...) API
X, Y, U, V, C, units, angles
```

```
pie(x,...) API
Z, explode, labels, colors, radius
```

```
text(x,y,text,...) API
x, y, text, va, ha, size, weight, transform
```

```
fill_between(x,...) API
X, Y1, Y2, color, where
```

## Advanced plots

```
step(X,Y,[fmt],...) API
X, Y, fmt, color, marker, where
```

```
boxplot(X,...) API
X, notch, sym, bootstrap, widths
```

```
errorbar(X,Y,xerr,yerr,...) API
X, Y, xerr, yerr, fmt
```

```
hist(X, bins, ...) API
X, bins, range, density, weights
```

```
violinplot(D,...) API
D, positions, widths, vert
```

```
barbs([X],[Y], U, V, ...) API
X, Y, U, V, C, length, pivot, sizes
```

```
eventplot(positions,...) API
positions, orientation, lineoffsets
```

```
hexbin(X,Y,C,...) API
X, Y, C, gridsize, bins
```

## Scales

```
ax.set_[xy]scale(scale,...) API
linear any values
log values > 0
symlog any values
logit 0 < values < 1
```

## Projections

```
subplot(..,projection=p) API
p='polar'
p='3d'
p=Orthographic()
from cartopy.crs import Cartographic
```

## Lines

```
linestyle or ls API
"-", "--", "-.", "-.", "(8,(0.61,2))
```

```
capstyle or dash_capstyle API
"butt", "round", "projecting"
```

## Markers

o	□	+	×	*	△	◇	▷	◁	▽	▽
◦	○	◆	▲	★	◆	◇	▷	◁	▽	▽
‘1’	‘2’	‘3’	‘4’	‘5’	‘6’	‘7’	‘8’	‘9’	‘10’	‘11’
‘S’	‘H’	‘C’	‘D’	‘P’	‘D’	‘P’	‘D’	‘P’	‘S’	‘H’
‘\$’	‘\$’	‘\$’	‘\$’	‘\$’	‘\$’	‘\$’	‘\$’	‘\$’	‘\$’	‘\$’

```
markervary API
10, [0,-1], (25,5), [0,25,-1]
```

## Colors

C0	C1	C2	C3	C4	C5	C6	C7	C8	C9	Cn
b	g	r	c	m	y	k	w	x	'x'	
Barkred	Firebrick	Crimson	Indianred	Salmon					'name'	
(1,0,0)	(1,0,0,0.75)	(1,0,0,0.5)	(1,0,0,0.25)						(R,G,B,[A])	
#FF0000	#FF0000BB	#FF00008B	#FF00008A						#RRGGBB[AA]	

```
'x' y' x'y' x''y'' x''y''
```

## Colormaps

```
plt.get_cmap(name) API
```

### Uniform

viridis	magma
plasma	

### Sequential

Greys	YlOrBr
	Wistia

### Diverging

Spectral	coolwarm
coolwarm	RdGy

### Qualitative

tab10	tab20

### Cyclic

twilight	

## Tick locators

```
from matplotlib import ticker
ax.[xy]axis.set_[minor|major]_locator(locator)
ticker.NullLocator()

ticker.MultipleLocator(0.5)
0.0 0.5 1.0 1.5 2.0 2.5 3.0 3.5 4.0 4.5 5.0
5 5.5 6 6.5 7 7.5 8 8.5 9 9.5 10 10.5 11 11.5 12 12.5 13 13.5 14 14.5 15 15.5 16 16.5 17 17.5 18 18.5 19 19.5 20 20.5 21 21.5 22 22.5 23 23.5 24 24.5 25 25.5 26 26.5 27 27.5 28 28.5 29 29.5 30 30.5 31 31.5 32 32.5 33 33.5 34 34.5 35 35.5 36 36.5 37 37.5 38 38.5 39 39.5 40 40.5 41 41.5 42 42.5 43 43.5 44 44.5 45 45.5 46 46.5 47 47.5 48 48.5 49 49.5 50

ticker.FixedLocator([0, 1, 5])
0 1 5 5.5 6 6.5 7 7.5 8 8.5 9 9.5 10 10.5 11 11.5 12 12.5 13 13.5 14 14.5 15 15.5 16 16.5 17 17.5 18 18.5 19 19.5 20 20.5 21 21.5 22 22.5 23 23.5 24 24.5 25 25.5 26 26.5 27 27.5 28 28.5 29 29.5 30 30.5 31 31.5 32 32.5 33 33.5 34 34.5 35 35.5 36 36.5 37 37.5 38 38.5 39 39.5 40 40.5 41 41.5 42 42.5 43 43.5 44 44.5 45 45.5 46 46.5 47 47.5 48 48.5 49 49.5 50

ticker.IndexLocator(base=0.5, offset=0.25)
0.25 0.5 0.75 1 1.25 1.5 1.75 2 2.25 2.5 2.75 3 3.25 3.5 3.75 4 4.25 4.5 4.75 5 5.25 5.5 5.75 6 6.25 6.5 6.75 7 7.25 7.5 7.75 8 8.25 8.5 8.75 9 9.25 9.5 9.75 10 10.25 10.5 10.75 11 11.25 11.5 11.75 12 12.25 12.5 12.75 13 13.25 13.5 13.75 14 14.25 14.5 14.75 15 15.25 15.5 15.75 16 16.25 16.5 16.75 17 17.25 17.5 17.75 18 18.25 18.5 18.75 19 19.25 19.5 19.75 20 20.25 20.5 20.75 21 21.25 21.5 21.75 22 22.25 22.5 22.75 23 23.25 23.5 23.75 24 24.25 24.5 24.75 25 25.25 25.5 25.75 26 26.25 26.5 26.75 27 27.25 27.5 27.75 28 28.25 28.5 28.75 29 29.25 29.5 29.75 30 30.25 30.5 30.75 31 31.25 31.5 31.75 32 32.25 32.5 32.75 33 33.25 33.5 33.75 34 34.25 34.5 34.75 35 35.25 35.5 35.75 36 36.25 36.5 36.75 37 37.25 37.5 37.75 38 38.25 38.5 38.75 39 39.25 39.5 39.75 40 40.25 40.5 40.75 41 41.25 41.5 41.75 42 42.25 42.5 42.75 43 43.25 43.5 43.75 44 44.25 44.5 44.75 45 45.25 45.5 45.75 46 46.25 46.5 46.75 47 47.25 47.5 47.75 48 48.25 48.5 48.75 49 49.25 49.5 49.75 50

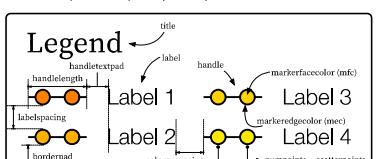
ticker.AutoLocator()
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50

ticker.MaxNLocator(n=4)
0.0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50

ticker.LogLocator(base=10, numticks=15)
10^1 10^2 10^3 10^4 10^5 10^6 10^7 10^8 10^9 10^10 10^11 10^12 10^13 10^14 10^15 10^16 10^17 10^18 10^19 10^20 10^21 10^22 10^23 10^24 10^25 10^26 10^27 10^28 10^29 10^30 10^31 10^32 10^33 10^34 10^35 10^36 10^37 10^38 10^39 10^40 10^41 10^42 10^43 10^44 10^45 10^46 10^47 10^48 10^49 10^50
```

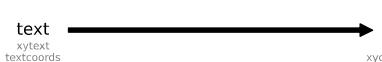
## Ornaments

```
ax.legend(...) API
handles, labels, loc, title, frameon
```



```
ax.colorbar(...) API
mappable, ax, cax, orientation
```

```
ax.annotate(...) API
text, xy, xytext, xycoords, textcoords, arrowprops
```



## Event handling

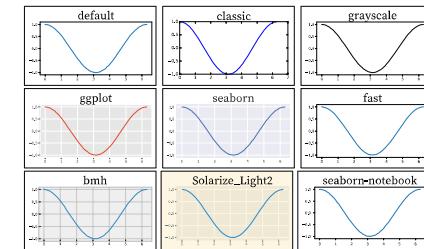
```
fig, ax = plt.subplots()
def on_click(event):
    print(event)
fig.canvas.mpl_connect('button_press_event', on_click)
```

## Animation

```
import matplotlib.animation as mpl
T = np.linspace(0, 2*np.pi, 100)
S = np.sin(T)
line, = plt.plot(T, S)
def animate(i):
    line.set_ydata(np.sin(T+i/50))
anim = mpl.FuncAnimation(
    plt.gcf(), animate, interval=5)
plt.show()
```

## Styles

```
plt.style.use(style)
```



## Quick reminder

```
ax.grid() API
ax.set_[xy]lim(vmin, vmax)
ax.set_[xy]label(label)
ax.set_[xy]ticks(ticks, [labels])
ax.set_[xy]ticklabels(labels)
ax.set_title(title)
ax.tick_params(width=10, ...)
ax.set_axis_[on|off]()
```

```
fig.suptitle(title)
fig.tight_layout()
plt.gcf(), plt.gca()
mpl.rc('axes', linewidth=1, ...)
[fig|ax].patch.set_alpha(0)
text=r'$\frac{e^{i\pi}}{2^n}$',
```

## Keyboard shortcuts

ctrl+s	Save	ctrl+w	Close plot
r	Reset view	f	Fullscreen 0/1
f	View forward	b	View back
p	Pan view	o	Zoom to rect
x	X pan/zoom	y	Y pan/zoom
g	Minor grid 0/1	G	Major grid 0/1
I	X axis log/linear	L	Y axis log/linear

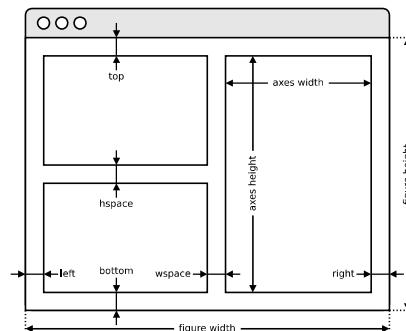
## Ten simple rules

1. Know Your Audience
2. Identify Your Message
3. Adapt the Figure
4. Captions Are Not Optional
5. Do Not Trust the Defaults
6. Use Color Effectively
7. Do Not Mislead the Reader
8. Avoid "Chartjunk"
9. Message Trumps Beauty
10. Get the Right Tool

## Axes adjustments

API

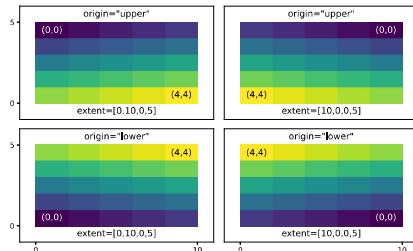
```
plt.subplots_adjust(...)
```



## Extent & origin

API

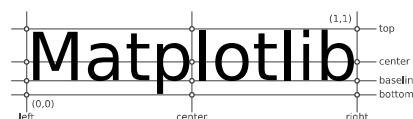
```
ax.imshow(extent=..., origin=...)
```



## Text alignments

API

```
ax.text(..., ha=..., va=..., ...)
```



## Text parameters

API

```
ax.text(..., family=..., size=..., weight=...)  
ax.text(..., fontproperties=...)
```

The quick brown fox

xx-large (1.73)

The quick brown fox

x-large (1.44)

The quick brown fox

large (1.20)

The quick brown fox

medium (1.00)

The quick brown fox

small (0.83)

The quick brown fox

x-small (0.69)

The quick brown fox

xx-small (0.58)

The quick brown fox jumps over the lazy dog

black (900)

The quick brown fox jumps over the lazy dog

bold (700)

The quick brown fox jumps over the lazy dog

semibold (600)

The quick brown fox jumps over the lazy dog

normal (400)

The quick brown fox jumps over the lazy dog

ultralight (100)

The quick brown fox jumps over the lazy dog

monospace

The quick brown fox jumps over the lazy dog

serif

The quick brown fox jumps over the lazy dog

sans

The quick brown fox jumps over the lazy dog

cursive

The quick brown fox jumps over the lazy dog

italic

The quick brown fox jumps over the lazy dog

normal

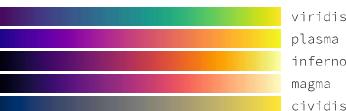
The quick brown fox jumps over the lazy dog

small-caps

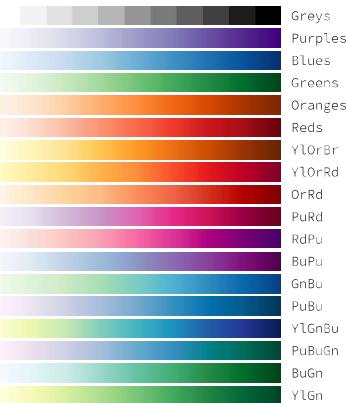
The quick brown fox jumps over the lazy dog

normal

## Uniform colormaps



## Sequential colormaps

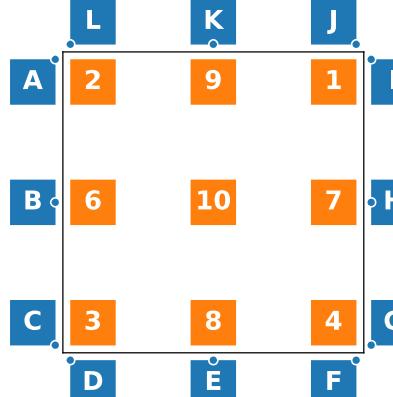


## Color names

API



## Legend placement



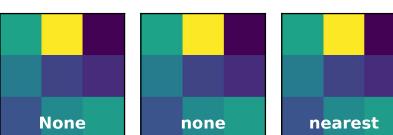
```
ax.legend(loc="string", bbox_to_anchor=(x,y))
```

2: upper left      9: upper center      1: upper right  
6: center left      10: center      7: center right  
3: lower left      8: lower center      4: lower right

A: upper right / (-0.1, 0.9)      B: center right / (-0.1, 0.5)  
C: lower right / (-0.1, 0.1)      D: upper left / (0.1, -0.1)  
E: upper center / (0.5, -0.1)      F: upper right / (0.9, -0.1)  
G: lower left / (1.1, 0.1)      H: center left / (1.1, 0.5)  
I: upper left / (1.1, 0.9)      J: lower right / (0.9, 1.1)  
K: lower center / (0.5, 1.1)      L: lower left / (0.1, 1.1)

## Image interpolation

API



## How do I ...

- ... resize a figure?  
→ `fig.set_size_inches(w, h)`
- ... save a figure?  
→ `fig.savefig("figure.pdf")`
- ... save a transparent figure?  
→ `fig.savefig("figure.pdf", transparent=True)`
- ... clear a figure/an axes?  
→ `fig.clear()` → `ax.clear()`
- ... close all figures?  
→ `plt.close("all")`
- ... remove ticks?  
→ `ax.set_[xy]ticks([])`
- ... remove tick labels?  
→ `ax.set_[xy]ticklabels([])`
- ... rotate tick labels?  
→ `ax.set_[xy]ticks(rotation=90)`
- ... hide top spine?  
→ `ax.spines['top'].set_visible(False)`
- ... hide legend border?  
→ `ax.legend(frameon=False)`
- ... show error as shaded region?  
→ `ax.fill_between(X, Y+error, Y-error)`
- ... draw a rectangle?  
→ `ax.add_patch(pt.Rectangle((0, 0), 1, 1))`
- ... draw a vertical line?  
→ `ax.axvline(x=0.5)`
- ... draw outside frame?  
→ `ax.plot(..., clip_on=False)`
- ... use transparency?  
→ `ax.plot(..., alpha=0.25)`
- ... convert an RGB image into a gray image?  
→ `gray = 0.2989*R + 0.5870*G + 0.1140*B`
- ... set figure background color?  
→ `fig.patch.set_facecolor("grey")`
- ... get a reversed colormap?  
→ `plt.get_cmap("viridis_r")`
- ... get a discrete colormap?  
→ `plt.get_cmap("viridis", 10)`
- ... show a figure for one second?  
→ `fig.show(block=False), time.sleep(1)`

## Performance tips

<code>scatter(X, Y)</code>	slow
<code>plot(X, Y, marker="o", ls="")</code>	fast
<code>for i in range(n): plot(X[i], Y[i])</code>	slow
<code>plot(sum([x[None] for x in X]), Y)</code>	fast
<code>cla(), imshow(...), canvas.draw()</code>	slow
<code>im.set_data(...), canvas.draw()</code>	fast

## Beyond Matplotlib

Seaborn: Statistical Data Visualization  
Cartopy: Geospatial Data Processing  
yt: Volumetric data Visualization  
mpld3: Bringing Matplotlib to the browser  
Databshader: Large data processing pipeline  
plotnine: A Grammar of Graphics for Python

Matplotlib Cheatsheets  
Copyright (c) 2021 Matplotlib Development Team  
Released under a CC-BY 4.0 International License

**NUMFOCUS**  
OPEN CODE = BETTER SCIENCE



# 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', 'F', 'F', 'M', '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 algs
```

### > 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=0.8).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.neighbors import KNeighborsClassifier
>>> 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(1,5),
...            "metric": ["euclidean", "cityblock"]}
>>> grid = GridSearchCV(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(1,6), "weights": ["uniform", "distance"]}
>>> search = RandomizedSearchCV(knn, param_distributions=params,
...                               cv=4, n_iter=8, random_state=5)
>>> search.fit(X_train, y_train)
>>> print(search.best_score_)
```



Learn Data Skills Online at [www.DataCamp.com](http://www.DataCamp.com)

# 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

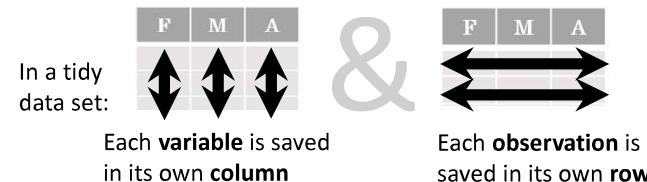
```
df = pd.DataFrame(
    {"a": [4, 5, 6],
     "b": [7, 8, 9],
     "c": [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [('d', 1), ('d', 2),
         ('e', 2)], names=['n', 'v']))
Create DataFrame with a MultiIndex
```

## Method Chaining

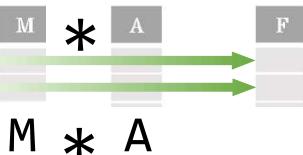
Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

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

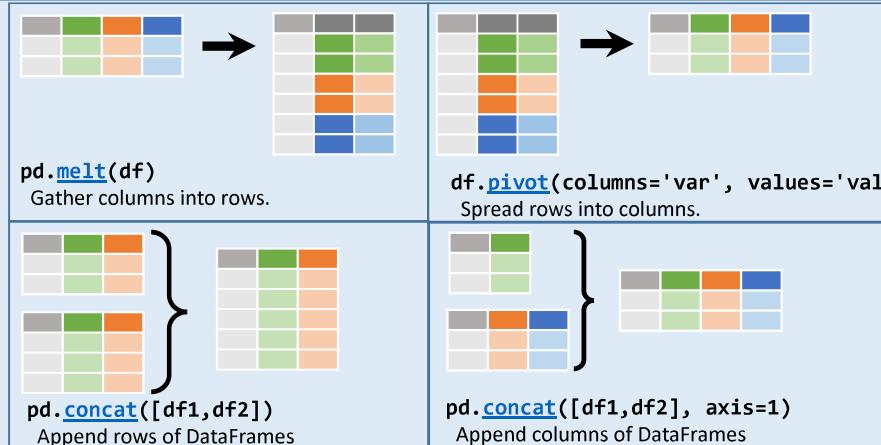
## Tidy Data – A foundation for wrangling in pandas



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

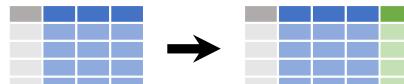
`df.dropna()`

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.

`min(axis=1)`

Element-wise min.

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

Trim values at input thresholds

`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.

## 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
C	3	NaN
D	NaN	T

`pd.merge(adf, bdf, how='outer', on='x1')`

Join data. Retain all values, all rows.

x1	x2
A	1
B	2

`adf[adf.x1.isin(bdf.x1)]`

All rows in adf that have a match in bdf.

x1	x2
C	3

`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
C	3
D	4



x1	x2
B	2
C	3
D	4

`pd.merge(ydf, zdf)`

Rows that appear in both ydf and zdf (Intersection).

x1	x2
A	1
C	3

`pd.merge(ydf, zdf, how='outer')`

Rows that appear in either or both ydf and zdf (Union).

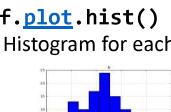
x1	x2
A	1

`pd.merge(ydf, zdf, how='outer', indicator=True)`

.query('\_merge == "left\_only")

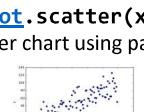
.drop(columns=['\_merge'])

Rows that appear in ydf but not zdf (Setdiff).



`df.plot.hist()`

Histogram for each column



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

Scatter chart using pairs of points