

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.randint(1000,10000)
>>> labels = np.random.randint(2,size=(1000,1))
>>> model = Sequential()
>>> model.add(Dense(32,
>                  activation='relu',
>                  input_dim=1000))
>>> 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,
>>> minst,
>>> cifar10,
>>> imbd
>>> (x_train,y_train), (x_test,y_test) = minst.load_data()
>>> (x_train,y_train2), (x_test2,y_test2) = boston_housing.load_data()
>>> (x_train,y_train3), (x_test3,y_test3) = cifar10.load_data()
>>> (x_train,y_train4), (x_test4,y_test4) = imbd.load_data(nm_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:-1]
>>> y = data[:, -1]
```

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_train3 = to_categorical(y_train3, num_classes)
>>> Y_test3 = to_categorical(y_test3, num_classes)
```

Model Architecture

Sequential Model

```
>>> from keras.models import Sequential
>>> model1 = 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),padding='same',input_shape=x_train.shape[1:]))
>>> 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_train5,y_train5,y_test5 = train_test_split(X,
>>> y,
>>> test_size=0.33,
>>> random_state=42)
```

Standardization/Normalization

```
>>> from sklearn.preprocessing import StandardScaler
>>> scaler = StandardScaler().fit(x_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
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=10,
>>> 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.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=10,
>>> validation_data=(x_test4,y_test4),
>>> callbacks=[early_stopping_monitor])
```

DataCamp

Learn Python for Data Science interactively



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



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), [(2,1), (4,5,6)]], dtype = float)
```

Initial Placeholders

```
>>> np.zeros((3,4))
Create an array of zeros
>>> np.ones((2,3,4),dtype=np.int16)
Create an array of ones
>>> np.empty((2,3,4))
Create an array of evenly spaced values (step value)
>>> np.linspace(0,2,9)
Create an array of evenly spaced values (number of samples)
>>> e = np.full((2,2),7)
Create a constant array
>>> f = np.eye(2)
Create a 2x2 identity matrix
>>> np.random.random((2,2))
Create an array with random values
>>> np.empty((3,2))
Create an empty array
```

I/O

Saving & Loading On Disk

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

Saving & Loading Text Files

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

Data Types

```
>>> np.int64
Signed 64-bit integer types
>>> np.float32
Double double-precision floating point
Complex numbers represented by 128 floats
>>> np.bool_
Boolean type storing TRUE and FALSE values
>>> np.object_
Python object type
>>> np.string_
Fixed-length string type
>>> np.unicode_
Fixed-length unicode type
```

Inspecting Your Array

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

Asking For Help

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

Array Mathematics

```
>>> a = b
array([[-0.5,  0. ,  0. ,  ..., 0. ,  ..., 0. ]])
>>> a + b
array([[ 2.5,  4. ,  6. ],
       [ 5. ,  7. ,  9. ]])
>>> np.add(a,b)
>>> a / b
array([[ 0.66666667,  1.        ,  1.        ,  ..., 1.        ,  ..., 1.        ]])
>>> np.divide(a,b)
>>> a * b
array([[ 1.5,  4. ,  9. ],
       [ 10. ,  18. ]])
>>> np.multiply(a,b)
>>> np.exp(b)
>>> np.sqrt(b)
>>> np.sin(a)
Print sines of an array
>>> np.cos(b)
>>> np.log(a)
>>> a.e.dot(f)
Element-wise dot product
>>> a * b
array([[ 1.5,  4. ,  9. ],
       [ 10. ,  18. ]])
>>> np.multiply(a,b)
>>> np.exp(b)
>>> np.sqrt(b)
>>> np.sin(a)
Print sines of an array
>>> np.cos(b)
>>> np.log(a)
>>> a.e.dot(f)
Element-wise dot product
```

Comparison

```
>>> a == b
Element-wise comparison
>>> a < 2
Element-wise comparison
>>> np.array_equal(a, b)
Element-wise comparison
```

Aggregate Functions

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

Copying Arrays

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

Sorting Arrays

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

Also see Lists

Select the element at the 2nd index

Select the element at row 0 column 2 (equivalent to b[1][2])

Select items at index 0 and 1

Select items at row 0 (equivalent to b[0,:])

Same as [1,:,:]

Reversed array

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
>>> a = np.transpose(b)
>>> i.T
Permute array dimensions
Permute array dimensions
```

Flatten the array

Reshape, but don't change data

Return a new array with shape (2,6)

Append items to an array

Insert items in an array

Delete items from an array

Concatenate arrays

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

DataCamp

Learn Python for Data Science interactively



Data Wrangling

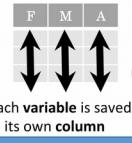
with pandas

Cheat Sheet

<http://pandas.pydata.org>

Tidy Data – A foundation for wrangling in pandas

In a tidy data set:



Each variable is saved in its own column



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



$M * A$

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	a	b	c
d	1	4	7	10
d	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

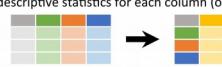
Method Chaining

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'))
```

Summarize Data

```
df['w'].value_counts()
Count number of rows with each unique value of variable
len(df)
# of rows in DataFrame.
df['w'].nunique()
# 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()	min()
Sum values of each object.	Minimum value in each object.
count()	max()
Count non-NA/null values of each object.	Maximum value in each object.
median()	mean()
Median value of each object.	Mean value of each object.
quantile([0.25, 0.75])	var()
Quantiles of each object.	Variance of each object.
apply(function)	std()
Apply function to each object.	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

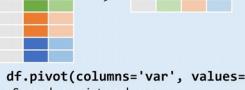
```
df.expanding()
Return an Expanding object allowing summary functions to be
applied cumulatively.
df.rolling(n)
Return a Rolling object allowing summary functions to be
applied to windows of length n.
```

Reshaping Data – Change the layout of a data set



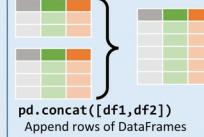
pd.melt(df)

Gather columns into rows.



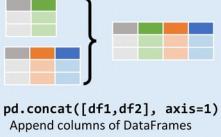
df.pivot(columns='var', values='val')

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(['Length','Height'], axis=1)

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
'^!(?i: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 .

Handling Missing Data

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)

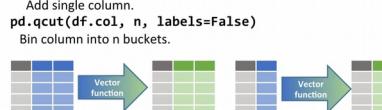
Compute and append one or more new columns.

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)

Absolute value.

min(axis=1)

Element-wise min.

abs()

Trim values at input thresholds

clip(lower=-10,upper=10)

Absolute value.

max()

Element-wise max.

min()

Element-wise min.

cumprod()

Cumulative product.

cumsum()

Cumulative sum.

cummax()

Cumulative max.

cummin()

Cumulative min.

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.

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.

Combine Data Sets

adf + bdf

+

=

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.

pd.merge(ydf, zdf,

how='outer')

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

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).

df.plot.hist()

Histogram for each column

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

Scatter chart using pairs of points

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.

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.

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 the summary functions listed above can also be applied to a group.

Additional GroupBy functions:

size()

Size of each group.

agg(function)

Aggregate group using function.

Windows

df.plot.hist()

Histogram for each column

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

Scatter chart using pairs of points

Windows

df.expanding()

Return an Expanding object allowing summary functions to be applied cumulatively.

df.rolling(n)

</div

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.



Use the following import convention:

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

Pandas Data Structures

Series

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



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

DataFrame

Columns

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

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

```
>>> 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)
```

```
>>> pd.to_csv('myDataFrame.csv')
```

Read and Write to Excel

```
>>> pd.read_excel('file.xlsx')
```

```
>>> pd.to_excel('dir/myDataFrame.xlsx', sheet_name='Sheet1')
```

Read multiple sheets from the same file

```
>>> xls = pd.ExcelFile('file.xls')
```

```
>>> df = pd.read_excel(xls, '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'

Select single value by row & column

By Label/Position

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

Country Brazil

Capital Brasilia

Population 207847528

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

0 Brussels

1 New Delhi

2 Brasilia

Select single row of subset of rows

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

'New Delhi'

Select a single column of subset of columns

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

'New Delhi'

Select rows and columns

Boolean Indexing

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

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

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

Series s where value is not > 1

s where value is <-1 or > 2

Use filter to adjust DataFrame

Setting

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

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(by='Country')
```

Sort by row or column index

```
>>> s.order()
```

Sort a series by its values

```
>>> 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.idmin() / df.idmax()
```

Minimum/Maximum index value

```
>>> df.describe()
```

Summary statistics

```
>>> df.mean()
```

Mean of values

```
>>> df.median()
```

Median of values

Applying Functions

```
>>> f = lambda x: x*x
```

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'])
```

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

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)
```

DataCamp

Learn Python for Data Science [Interactively](#)



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+5),2,3])
```

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

Index Tricks

```
>>> np.mgrid[0:5,0:5]
```

Create a dense meshgrid

```
>>> np.ogrid[0:2,0:2]
```

Create an open meshgrid

```
>>> np.r_[3,0]*5,-1:1:10j]
```

Stack arrays vertically (row-wise)

```
>>> np.c_[b,c]
```

Stack arrays horizontally (column-wise)

```
>>> np.vsplit(d,2)
```

Split the array horizontally at the 2nd index

```
>>> np.hsplit(c,2)
```

Split the array vertically at the 2nd index

```
>>> np.transpose(b)
```

Permute array dimensions

```
>>> b.flatten()
```

Flatten the array

```
>>> np.vstack((a,b))
```

Stack arrays vertically (column-wise)

```
>>> np.vstack((a,b))
```

Stack arrays vertically (row-wise)

```
>>> np.hsplit(c,2)
```

Split the array vertically (row-wise)

```
>>> np.vsplit(d,2)
```

Split the array horizontally (column-wise)

```
>>> np.concatenate((a,b))
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1)
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=0)
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=-1)
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='C')
```

Concatenate arrays along a given axis

```
>>> np.concatenate((a,b),axis=1,order='F')
```

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.mgrid[-3:3:100j, -3:3:100j]  
>>> V = 1 + X**2 + 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.rcParams['figure.figsize'])
```

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()  
>>> ax = 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)  
>>> ax.vlines([0,0.5,1], [3,4,5])  
>>> ax.patches([0,0.5,1], [0,1,2])  
>>> axes[1,1].axhline(0.45)  
>>> axes[0,1].axvline(0.65)  
>>> ax.fill(x,y,color='blue')  
>>> ax.fill_between(x,y,color='yellow')
```

Draw points with lines or markers connecting them
Draw unconnected points, scaled or colored
Plot horizontal rectangles (constant width)
Draw a horizontal line across axes
Draw a vertical line across axes
Draw filled polygons
Fill between y-values and o

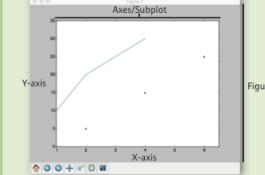
2D Data or Images

```
>>> fig, ax = plt.subplots()  
>>> im = ax.imshow(img, interpolation='bicubic',  
    cmap='gist_earth',  
    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] Step 1  
>>> fig = plt.figure() Step 2  
>>> ax = fig.add_subplot(111) Step 3  
>>> ax.plot(x, y, color='lightblue', linewidth=3) Step 3  
>>> ax.scatter([2, 3, 4], [15, 15, 25],  
    color='darkgreen',  
    marker='^') Step 4  
>>> ax.set_xlim(1, 6.5) Step 4  
>>> plt.savefig('foo.png') Step 5  
>>> plt.show() Step 6
```

4 Customize Plot

Colors, Color Bars & Color Maps

```
>>> plt.plot(x, y, alpha=0.4)  
>>> ax.set_alpha(0.5)  
>>> fig.colorbar(im, orientation='horizontal')  
>>> im = ax.imshow(img,  
    cmap='seismic')
```

Markers

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

Linestyles

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

Text & Annotations

```
>>> ax.text(1,-2.1,  
    "Simple Graph",  
    style='italic')  
>>> ax.annotate("Sine",  
    xy=(8, 0),  
    xytext=(10.5, 0),  
    textcoords="data",  
    arrowprops=dict(arrowstyle=">-",  
    connectionstyle="arc3"))
```

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.bxpplot(y)  
>>> ax3.violinplot(z)
```

Plot a histogram
Make a box and whisker plot
Make a violin plot

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))
```

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

5 Plot

Save Plot

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

6 Show Plot

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

Close & Clear

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

DataCamp
Learn Python for Data Science interactively



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.cross_validation import train_test_split  
>>> from sklearn.metrics import accuracy_score  
>>> iris = datasets.load_iris()  
>>> X_train, X_test, y_train, y_test = train_test_split(  
    X, y, test_size=0.4, 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 DataFrames, are also acceptable.

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

Training And Test Data

```
>>> from sklearn.cross_validation 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  
>>> normalizer = Normalizer().fit(X_train)  
>>> normalized_X = normalizer.transform(X_train)  
>>> normalized_X_test = normalizer.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
```

KNN

```
>>> from sklearn.neighbors import KNeighborsClassifier
```

```
>>> knn = 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.fit_transform(X_train)
```

Fit the model to the data

Fit the model to the data

Fit the model to the data

Fit to data, then transform it

Prediction

Supervised Estimators

```
>>> y_pred = lr.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 algs

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=np.nan, strategy='mean', axis=0)
```

```
>>> imp.fit_transform(X_train)
```

Generating Polynomial Features

```
>>> from sklearn.preprocessing import PolynomialFeatures
```

```
>>> poly = PolynomialFeatures(5)
```

```
>>> poly.fit_transform(X)
```

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)
```

Metric scoring functions

Classification Report

```
>>> from sklearn.metrics import classification_report
```

```
>>> print(classification_report(y_test, y_pred))
```

Precision, recall, f-score and support

Confusion Matrix

```
>>> from sklearn.metrics import confusion_matrix
```

```
>>> print(confusion_matrix(y_test, y_pred))
```

Regression Metrics

Mean Absolute Error

```
>>> from sklearn.metrics import mean_absolute_error
```

```
>>> y_true = [3, -0.5, 2]
```

```
>>> mean_absolute_error(y_true, y_pred)
```

Mean Squared Error

```
>>> from sklearn.metrics import mean_squared_error
```

```
>>> mean_squared_error(y_test, y_pred)
```

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 Score

```
>>> 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"]}
```

```
>>> research = RandomizedSearchCV(estimator=knn,  
    param_distributions=params,
```

```
>>> n_iter=8,  
    random_state=5)
```

```
>>> research.fit(X_train, y_train)
```

```
>>> print(research.best_score_)
```

DataCamp

Learn Python for Data Science interactively

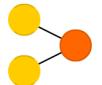


A mostly complete chart of
Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool

Perceptron (P)



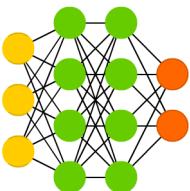
Feed Forward (FF)



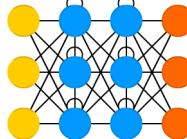
Radial Basis Network (RBF)



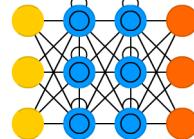
Deep Feed Forward (DFF)



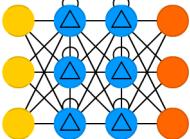
Recurrent Neural Network (RNN)



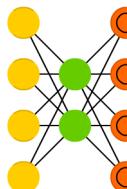
Long / Short Term Memory (LSTM)



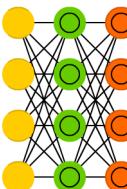
Gated Recurrent Unit (GRU)



Auto Encoder (AE)



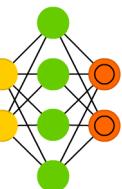
Variational AE (VAE)



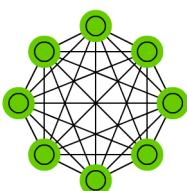
Denoising AE (DAE)



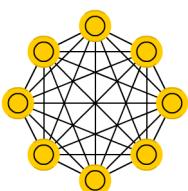
Sparse AE (SAE)



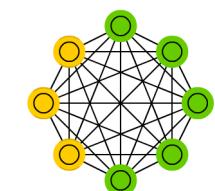
Markov Chain (MC)



Hopfield Network (HN)



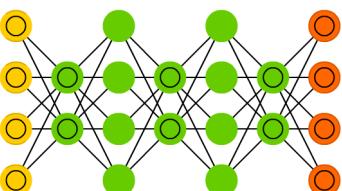
Boltzmann Machine (BM)



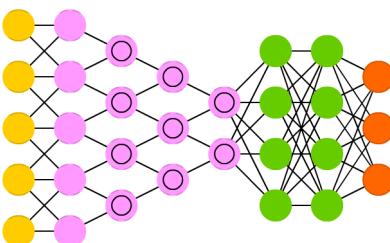
Restricted BM (RBM)



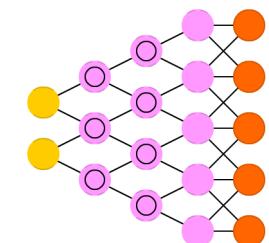
Deep Belief Network (DBN)



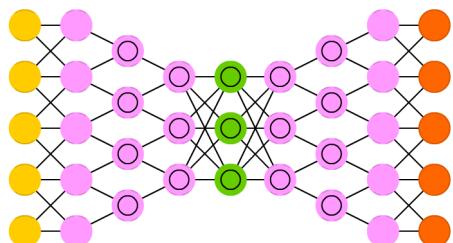
Deep Convolutional Network (DCN)



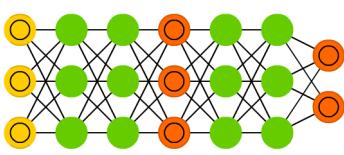
Deconvolutional Network (DN)



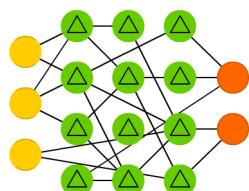
Deep Convolutional Inverse Graphics Network (DCIGN)



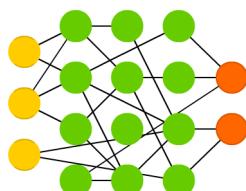
Generative Adversarial Network (GAN)



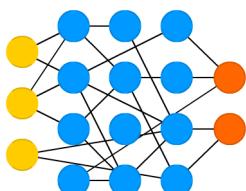
Liquid State Machine (LSM)



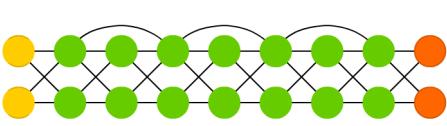
Extreme Learning Machine (ELM)



Echo State Network (ESN)



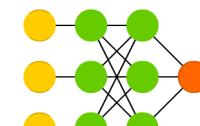
Deep Residual Network (DRN)



Kohonen Network (KN)



Support Vector Machine (SVM)



Neural Turing Machine (NTM)



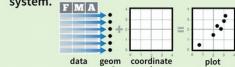
Data Visualization with ggplot2

Cheat Sheet

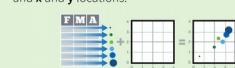


Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same few components: a **data set**, a set of **geoms**—visual marks that represent data points, and a **coordinate system**.



To display data values, map variables in the data set to aesthetic properties of the geom like **size**, **color**, and **x** and **y** locations.



Build a graph with **plot()** or **ggplot()**

```
aesthetic mappings   data   geom
ggplot(=cty, =hwy, color = cyl, data = mpg, geom = "point")
Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.
```

ggplot(data = mpg, aes(x ~ cyl, y = hwy))

Begins a plot that you finish by adding layers to. No defaults, but provides more control than **plot()**.

```
data
ggplot(mpg, aes(hwy, cyl)) +
  geom_point(aes(color = cyl)) +
  geom_smooth(method = "lm") +
  coord_cartesian() +
  scale_color_gradient() +
  theme_bw()
  add layers, elements with +
  layer specific mappings
  default stat
  additional elements
```

Add a new layer to a plot with a **geom_*** or **stat_*** function. Each provides a geom, a set of aesthetic mappings, and a default stat and position adjustment.

last_plot()

Returns the last plot

ggsave("plot.png", width = 5, height = 5)

Saves last plot as 5 x 5" file named "plot.png" in working directory. Matches file type to file extension.

RStudio® is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 844-448-1212 • studio.com

Geoms - Use a geom to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

One Variable

Continuous

```
a <- ggplot(mpg, aes(hwy))
a + geom_area(stat = "bin")
x, y, alpha, color, fill, linetype, size
a + geom_density(kernel = "gaussian")
x, y, alpha, color, fill, linetype, size, weight
a + geom_dotplot()
x, y, alpha, color, fill
a + geom_freqpoly()
x, y, alpha, color, linetype, size
a + geom_freqpoly(aes(..density..))
a + geom_histogram(binwidth = 5)
x, y, alpha, color, fill, linetype, size, weight
a + geom_histogram(aes(..density..))
  
```

Discrete

```
b <- ggplot(mpg, aes(fl))
b + geom_bar()
x, alpha, color, fill, linetype, size, weight
  
```

Graphical Primitives

```
c <- ggplot(map, aes(long, lat))
c + geom_polygon(aes(group = group))
x, y, alpha, color, fill, linetype, size
  
```

```
d <- ggplot(economics, aes(date, unemploy))
d + geom_path(lineend = "butt",
linejoin = "round", linemtire = 1)
x, y, alpha, color, linetype, size
d + geom_ribbon(aes(ymin = unemploy - 900,
ymax = unemploy + 900))
x, ymax, ymin, alpha, color, fill, linetype, size
  
```

```
e <- ggplot(seals, aes(x = long, y = lat))
e + geom_segment(aes(
xend = long + delta_long,
yend = lat + delta_lat))
x, end, y, end, alpha, color, linetype, size
e + geom_rect(aes(xmin = long, ymin = lat,
xmax = long + delta_long,
ymax = lat + delta_lat))
xmax, xmin, ymax, ymin, alpha, color, fill,
linetype, size
  
```

```
f <- ggplot(diamonds, aes(cut, color))
f + geom_jitter()
x, y, alpha, color, fill, shape, size
  
```

Learn more at docs.ggplot2.org • ggplot2 0.9.3.1 • Updated: 3/15

Two Variables

Continuous X, Continuous Y

```
f <- ggplot(mpg, aes(cty, hwy))
f + geom_blank()
f + geom_jitter()
x, y, alpha, color, fill, shape, size
f + geom_point()
x, y, alpha, color, fill, shape, size
f + geom_quantile()
x, y, alpha, color, linetype, size, weight
f + geom_rug(sides = "bl")
alpha, color, linetype, size
f + geom_smooth(model = lm)
x, y, alpha, color, fill, linetype, size, weight
C f + geom_text(aes(label = cty))
x, y, label, alpha, angle, color, family, fontface,
hjust, lineheight, size, vjust
  
```

Continuous Bivariate Distribution

```
i + geom_bin2d(binwidth = c(5, 0.5))
xmax, xmin, ymax, ymin, alpha, color, fill,
linetype, size, weight
i + geom_density2d()
x, y, alpha, colour, linetype, size
i + geom_hex()
x, y, alpha, colour, fill size
  
```

Continuous Function

```
j <- ggplot(economics, aes(date, unemploy))
j + geom_area()
x, y, alpha, color, fill, linetype, size
j + geom_line()
x, y, alpha, color, linetype, size
j + geom_step(direction = "hv")
x, y, alpha, color, linetype, size
  
```

Visualizing error

```
df <- data.frame(grp = c("A", "B"), fit = 4.5; se = 1.2)
k <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
  
```

```
k + geom_crossbar(fatten = 2)
x, y, ymax, ymin, alpha, color, fill, linetype,
size
k + geom_errorbar()
x, ymax, ymin, alpha, color, linetype, size,
width (also geom_errorbarh())
k + geom_linerange()
x, ymin, ymax, alpha, color, linetype, size
k + geom_pointrange()
x, y, ymin, ymax, alpha, color, fill, linetype,
shape, size
  
```

Maps

```
data <- data.frame(murder = USArrests$Murder,
state = tolower(rownames(USArrests)))
map <- map_data("state")
l <- ggplot(data, aes(fill = murder))
l + geom_map(aes(map_id = state), map = map) +
expand_limits(x = map$long, y = map$lat)
map_id, alpha, color, fill, linetype, size
  
```

Three Variables

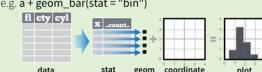
```
sealsSz <- with(seals, sqrt(delta_long^2 + delta_lat^2))
m <- ggplot(seals, aes(long, lat))
  
```

```
m + geom_contour(aes(z = z))
x, y, z, alpha, colour, linetype, size, weight
  
```

```
m + geom_raster(aes(fill = z), hjust = 0.5,
vjust = 0.5, interpolate = FALSE)
x, y, alpha, fill
m + geom_tile(aes(fill = z))
x, y, alpha, color, fill, linetype, size
  
```

Stats

Some plots visualize a **transformation** of the original data set. Use a **stat** to choose a common transformation to visualize, e.g. **a + geom_bar(stat = "bin")**



Each stat creates additional variables to map aesthetics to. These variables use a common **.name..** syntax.

stat functions and geom functions both combine a stat with a geom to make a layer, i.e. **stat_bin(geom="bar")** does the same as **geom_bar(stat="bin")**

```
stat function    layer specific mappings    variable created by transformation
  
```

```
i + stat_density2d(aes(fill = ..level..),
geom = "polygon", n = 100)
geom for layer    parameters for stat
  
```

```
a + stat_bin(binwidth = 1, origin = 10)
x, y, fill | count.., ..count., ..density..
a + stat_bin(binwidth = 1, bins = 2)
x, y, fill | count..
a + stat_density(adjust = 1, kernel = "gaussian")
x, y, ..scaled..
  
```

```
f + stat_bin2d(bins = 30, drop = TRUE)
x, y, fill | count.., density..
f + stat_hexbin(bins = 30)
x, y, fill | count..
f + stat_density2d(contour = TRUE, n = 100)
x, y, size | ..level..
  
```

```
m + stat_contour(aes(z = z))
x, y, z, order | level.
  
```

```
m + stat_spoke(aes(radius = z, angle = z))
angle, radius, x, xend, y, yend | x, ..xend., y, ..yend..
m + stat_summary_hex(aes(z = z), bins = 30, fun = mean)
x, y, z, fill | value.
  
```

```
m + stat_summary2d(aes(z = z), bins = 30, fun = mean)
x, y, z, fill | ..value..
  
```

```
g + stat_boxplot(coef = 1.5)
Comparisons
x, y | lower, ..middle., upper, ..outliers..
g + stat_ydensity(aes(ydensity = 1, kernel = "gaussian", scale = "area"))
x, y | ..density.., ..count.., ..n.., ..violinwidth.., ..width..
  
```

```
f + stat_ecdf(p = 40)
Functions
x, y | ..x.., ..y..
  
```

```
f + stat_quantile(quartiles = c(0.25, 0.5, 0.75), formula = y ~ log(x),
method = "rq")
x, y | quantile, ..x.., ..y..
  
```

```
f + stat_smooth(method = "auto", formula = y ~ x, se = TRUE, n = 80,
fullrange = FALSE, level = 0.95)
x, y | ..se.., ..x.., ..y.., ..ymean.., ..ymin..
  
```

```
ggplot() + stat_function(aes(x ~ .33))
General Purpose
fun = dnorm, n = 101, args = list(d = 0.5)
x | ..x..
  
```

```
f + stat_identity()
ggplot() + stat_qq(aes(sample = 1:100), distribution = qt,
dparams = list(fit = 0))
sample, x, y | ..x.., ..y..
  
```

```
f + stat_kde()
x, y | ..size..
  
```

```
f + stat_summary(fun.data = "mean_cl_boot")
f + stat_uniques()
  
```

Scales

Scales control how a plot maps data values to the visual values of an aesthetic. To change the mapping, add a custom scale.

```
n <- b + geom_bar(aes(fill = fill))
n + scale_fill_manual(values = c("skyblue", "royalblue", "blue", "navy"),
limits = c("D", "E", "P", "R"), breaks = c("D", "E", "P", "R"),
name = "fuel", labels = c("D", "E", "P", "R"))
  
```

```
range of values to include in mapping
title to use in legend/axis
labels to use in legend/axis
breaks to use in legend/axis
  
```

General Purpose scales

Use with any aesthetic: alpha, color, fill, linetype, shape, size

```
scale_*_continuous() - map cont values to visual values
scale_*_discrete() - map discrete values to visual values
scale_*_identity() - use data values as visual values
scale_*_manual() - map values to c()
  
```

```
manually chosen visual values
X and Y location scales
Use with x or y aesthetics (x shown here)
scale_x_date(aes(labels = date_format("%m/%d/%Y"),
breaks = date_breaks("2 weeks")) - treat x values as dates. See ?strptime for label formats.
scale_x_datetime() - treat x values as date times. Use same arguments as scale_x_date().
scale_x_log10() - Plot on log10 scale
scale_x_reverse() - Reverse direction of x axis
scale_x_sqrt() - Plot on square root scale
  
```

```
Color and fill scales
Discrete
  
```

```
a <- b + geom_bar(aes(fill = ..level..))
a + scale_fill_brewer(palette = "Set1")
  
```

```
For palette choices: library(RColorBrewer); display.brewer.all()
  
```

```
f + scale_fill_grey(start = 0.2, end = 0.8,
na.value = "#red")
  
```

```
g + scale_fill_hexagonal(colors = terrain.colors(6))
Also rainbow(), heat.colors(), topo.colors(), cm.colors(),
PCorBrewer.brewer.pal()
  
```

```
  
```

```
Shape scales
Manual shape values
  
```

```
p <- f + geom_point(aes(shape = fill))
p + scale_shape_manual(values = c(19, 21))
  
```

```
  
```

```
p + scale_shape_hexagonal()
  
```

```
  
```

```
  
```

```
Size scales
Manual size values
  
```

```
q <- f + geom_point(aes(size = cyl))
q + scale_size_area(max = 1000, aes(size = cyl))
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```

```
  
```