

Python For Data Science Cheat Sheet

Keras

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

Keras

Keras is a powerful and easy-to-use deep learning library for Theano and TensorFlow that provides a high-level neural networks API to develop and evaluate deep learning models.

A Basic Example

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

Data

Also see NumPy, Pandas & Scikit-Learn

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

Keras Data Sets

```
>>> from keras.datasets import boston_housing, mnist, cifar10, imdb
>>> (x_train,y_train),(x_test,y_test) = mnist.load_data()
>>> (x_train2,y_train2),(x_test2,y_test2) = boston_housing.load_data()
>>> (x_train3,y_train3),(x_test3,y_test3) = cifar10.load_data()
>>> (x_train4,y_train4),(x_test4,y_test4) = imdb.load_data(num_words=20000)
>>> num_classes = 10
```

Other

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

Preprocessing

Sequence Padding

```
>>> from keras.preprocessing import sequence
>>> x_train4 = sequence.pad_sequences(x_train4,maxlen=80)
>>> x_test4 = sequence.pad_sequences(x_test4,maxlen=80)
```

One-Hot Encoding

```
>>> from keras.utils import to_categorical
>>> Y_train = to_categorical(y_train, num_classes)
>>> Y_test = to_categorical(y_test, num_classes)
>>> Y_train3 = to_categorical(y_train3, num_classes)
>>> Y_test3 = to_categorical(y_test3, num_classes)
```

Model Architecture

Sequential Model

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

Multi-layer Perceptron (MLP)

Binary Classification

```
>>> from keras.layers import Dense
>>> model.add(Dense(12,
                  input_dim=8,
                  kernel_initializer='uniform',
                  activation='sigmoid'))
>>> 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=(784,1)))
>>> model.add(Dense(512, activation='relu'))
>>> model.add(Dropout(0.2))
>>> model.add(Dense(10, activation='softmax'))
```

Regression

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

Convolutional Neural Network (CNN)

```
>>> from keras.layers import Activation, Conv2D, MaxPooling2D, Flatten
>>> model2.add(Conv2D(32, (3,3), padding='same', input_shape=x_train.shape[1:]))
>>> model2.add(Activation('relu'))
>>> model2.add(Conv2D(32, (3,3)))
>>> model2.add(Activation('relu'))
>>> model2.add(MaxPooling2D(pool_size=(2,2)))
>>> model2.add(Dropout(0.25))
>>> model2.add(Conv2D(64, (3,3), padding='same'))
>>> model2.add(Activation('relu'))
>>> model2.add(Conv2D(64, (3,3)))
>>> model2.add(Activation('relu'))
>>> model2.add(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'))
```

Inspect Model

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

Model output shape
Model summary representation
Model configuration
List all weight tensors in the model

Compile Model

MLP: Binary Classification

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

MLP: Multi-Class Classification

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

MLP: Regression

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

Recurrent Neural Network

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

Model Training

```
>>> model3.fit(x_train4,
             y_train4,
             batch_size=32,
             epochs=15,
             verbose=1,
             validation_data=(x_test4,y_test4))
```

Evaluate Your Model's Performance

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

Prediction

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

Save/Reload Models

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

Model Fine-tuning

Optimization Parameters

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

Early Stopping

```
>>> from keras.callbacks import EarlyStopping
>>> early_stopping_monitor = EarlyStopping(patience=2)
>>> model3.fit(x_train4,
            y_train4,
            batch_size=32,
            epochs=15,
            validation_data=(x_test4,y_test4),
            callbacks=[early_stopping_monitor])
```

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

1D array



2D array



3D array



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))
```

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)
>>> b.ndim
>>> e.size
>>> b.dtype
>>> b.dtype.name
>>> b.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. ],
        [  1. , -1. , -1. ]])
>>> np.subtract(a,b)
>>> b + a
>>> array([[ 2.5,  4. ,  6. ],
        [  5. ,  7. ,  9. ]])
>>> np.add(b,a)
>>> a / b
>>> array([[ 0.6666667,  1. ,  1. ],
        [ 0.25,  0.4 ,  0.5 ]])
>>> np.divide(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.]])
```

Comparison

```
>>> a == b
>>> array([[True,  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.corcoef()
>>> 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

Also see Lists

Subsetting

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

Select the element at the 2nd index
Select the element at row 0 column 2 (equivalent to b[0][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][:, :])
Same as [1, :, :]
Reversed array a
Select elements from a less than 2

Slicing

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

Boolean Indexing

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

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.,  2.,  3., 1.5],
        [ 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

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. ]])
```

Concatenate arrays
Stack arrays vertically (row-wise)

Stacking Arrays

```
>>> np.r_[e,f]
>>> np.hstack((e,f))
>>> array([[ 7.,  7.,  1.,  0. ],
        [ 1.7,  7.,  0.,  1.1]])
```

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

Column-wise Arrays

```
>>> np.column_stack((a,d))
>>> array([[ 1, 10, ],
        [ 2, 15, ],
        [ 3, 20]])
```

Create stacked column-wise arrays

Splitting Arrays

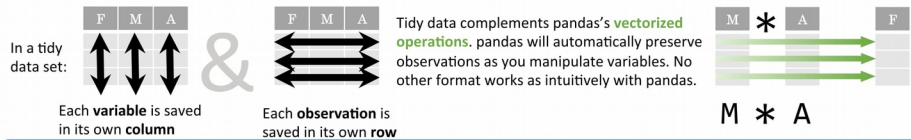
```
>>> np.hsplit(a,3)
>>> (array([1, 2, 3]), array([2]), array([3]))
>>> np.vsplit(c,2)
>>> array([[1.5, 2., 1. ],
        [ 4., 5., 6. ]]),
        array([[ 3., 2., 3. ],
        [ 4., 5., 6.]])
```

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

DataCamp

Learn Python for Data Science interactively

Tidy Data – A foundation for wrangling in pandas



Syntax – Creating DataFrames

```
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.

```
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']))
```

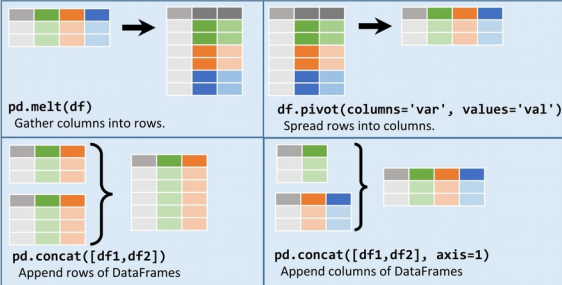
Create DataFrame with a MultiIndex

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'))
```

Reshaping Data – Change the layout of a data set



```
df.sort_values('mpg')
df.sort_values('mpg', ascending=False)
df.rename(columns = {'y': 'year'})
df.sort_index()
df.reset_index()
df.drop(['Length', 'Height'], axis=1)
```

Subset Observations (Rows)

```
df[df.Length > 7]
df.sample(frac=0.5)
df.sample(n=10)
df.drop_duplicates()
df.iloc[10:20]
df.head(n)
df.tail(n)
```

Subset Variables (Columns)

```
df[['width', 'length', 'species']]
df['width'] or df.width
df.filter(regex='regex')
```

regex (Regular Expressions)	Examples
.*	Matches strings containing a period '.'
.*Length\$	Matches strings ending with word 'Length'
.*Sepal1	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'

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

Summarize Data

```
df['w'].value_counts()
len(df)
df['w'].nunique()
df.describe()
```

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()
count()
median()
quantile([0.25, 0.75])
apply(function)
```

Group Data

```
df.groupby(by="col")
df.groupby(level="ind")
agg(function)
```

Handling Missing Data

```
df.dropna()
df.fillna(value)
```

Make New Columns

```
df.assign(Area=lambda df: df.Length*df.Height)
df['Volume'] = df.Length*df.Height*df.Depth
pd.qcut(df.col, n, labels=False)
```

```
max(axis=1)
clip(lower=-10, upper=10)
shift(1)
rank(method='dense')
rank(method='min')
rank(pct=True)
clip(lower=-10, upper=10)
abs()
```

```
min(axis=1)
cumsum()
cummax()
cummin()
cumprod()
```

Combine Data Sets

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

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

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

```
ydf + zdf
```

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

Windows

```
df.expanding()
df.rolling(n)
```

Plotting

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

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

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

Asking For Help

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

Selection

Also see NumPy Arrays

Getting

```
>>> s['b']
-5
>>> df[1:]
   Country  Capital  Population
1  India  New Delhi  1303171035
2  Brazil  Brasilia  207847528
```

Get one element

Get subset of a DataFrame

Selecting, Boolean Indexing & Setting

```
By Position
>>> df.iloc[0,0]
'Belgium'
>>> df.loc[0,0]
'Belgium'
```

Select single value by row & column

```
By Label
>>> df.loc[0, ['Country']]
'Belgium'
>>> df.at[0, ['Country']]
'Belgium'
```

Select single value by row & column labels

```
By Label/Position
>>> df.ix[2]
Country      Brazil
Capital      Brasilia
Population    207847528
```

Select single row of subset of rows

```
>>> df.ix[:, 'Capital']
0  Brussels
1  New Delhi
2  Brasilia
```

Select a single column of subset of columns

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

Select rows and columns

Boolean Indexing

```
>>> s[s > 1]
a    3
b   -5
c    7
d    4
>>> s[(s < -1) | (s > 2)]
a    3
c    7
>>> df[df['Population'] > 1200000000]
```

Series *s* where value is not > 1 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 index value
>>> df.idxmin()/df.idxmax()
Minimum/Maximum index value
>>> df.describe()
Summary statistics
>>> df.mean()
Mean of values
>>> df.median()
Median of values
```

Applying Functions

```
>>> f = lambda x: x*2
>>> df.apply(f)
Apply function
>>> df.applymap(f)
Apply function element-wise
```

Data Alignment

Internal Data Alignment

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

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
>>> s + s3
a    10.0
b     NaN
c     5.0
d     7.0
```

Arithmetic Operations with Fill Methods

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

```
>>> s.add(s3, fill_value=0)
a    10.0
b     -5.0
c     5.0
d     7.0
>>> s.sub(s3, fill_value=2)
a    12.0
b     -7.0
c     3.0
d     4.0
>>> s.div(s3, fill_value=4)
a    2.5
b     NaN
c    1.25
d    1.75
>>> s.mul(s3, fill_value=3)
a    21.0
b     NaN
c    15.0
d    21.0
```

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+5j, 2j, 3j), (4,5,6)])
>>> c = np.array([(1.5, 2, 3), (4, 5, 6)], [(3, 2, 1), (4, 5, 6)])
Index Tricks
>>> np.mgrid[0:5,0:5]
Create a dense meshgrid
>>> np.ogrid[0:2,0:2]
Create an open meshgrid
>>> np.c_[3, [0] * 5, -11:10j]
Stack arrays vertically (row-wise)
>>> np.c_[b,c]
Create stacked column-wise arrays
Shape Manipulation
>>> np.transpose(b)
Permute array dimensions
>>> b.flatten()
Flatten the array
>>> np.hstack((b,c))
Stack arrays horizontally (column-wise)
>>> np.vstack((a,b))
Stack arrays vertically (row-wise)
>>> np.hsplit(c,2)
Split the array horizontally at the 2nd index
>>> np.vsplit(d,2)
Split the array vertically at the 2nd index
Polynomials
>>> from numpy import polyval
>>> p = polyval([3,4,5])
Create a polynomial object
Vectorizing Functions
>>> def myfunc(a):
    if a < 0:
        return a*2
    else:
        return a/2
>>> np.vectorize(myfunc)
Vectorize functions
Type Handling
>>> np.real(b)
Return the real part of the array elements
>>> np.imag(b)
Return the imaginary part of the array elements
>>> np.real_if_close(c, tol=1000)
Return a real array if complex parts close to 0
>>> np.cast['F'](np.pi)
Cast object to a data type
Other Useful Functions
>>> np.angle(b, deg=True)
Return the angle of the complex argument
>>> g = np.linspace(0, np.pi, num=5)
Create an array of evenly spaced values (number of samples)
>>> g[3] += np.pi
>>> np.unwrap(g)
>>> np.linspace(0, 10, 3)
Create an array of evenly spaced values (log scale)
>>> np.select([c<4], [c*2])
Return values from a list of arrays depending on conditions
>>> misc.factorial(a)
Factorial
>>> misc.comb(10, 3, exact=True)
Combine N things taken at k time
>>> misc.central_diff_weights(3)
Weights for Np-point central derivative
>>> misc.derivative(myFunc, 1, 0)
Find the n-th derivative of a function at a point
```

Linear Algebra

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

Also see NumPy

```
>>> from scipy import linalg, sparse
```

Creating Matrices

```
>>> A = np.matrix(np.random.random((2,2)))
>>> B = np.asmatrix(b)
>>> C = np.mat(np.random.random((10,5)))
>>> D = np.mat([[3,4], [5,6]])
```

Matrix Functions

```
Addition
>>> np.add(A,D)
Addition
Subtraction
>>> np.subtract(A,D)
Subtraction
Division
>>> np.divide(A,D)
Division
Multiplication
>>> A @ D
Multiplication operator (Python)
>>> np.multiply(D,A)
Multiplication
>>> np.dot(A,D)
Dot product
>>> np.vdot(A,D)
Vector dot product
>>> np.inner(A,D)
Inner product
>>> np.outer(A,D)
Outer product
>>> np.tensordot(A,D)
Tensor dot product
>>> np.kron(A,D)
Kronecker product
Exponential Functions
>>> linalg.expm(A)
Matrix exponential
>>> linalg.exp2(A)
Matrix exponential (Taylor Series)
>>> linalg.expml(D)
Matrix exponential (eigenvalue decomposition)
Logarithm Function
>>> linalg.logm(A)
Matrix logarithm
Trigonometric Functions
>>> linalg.sinm(D)
Matrix sine
>>> linalg.cosm(D)
Matrix cosine
>>> linalg.tanm(A)
Matrix tangent
Hyperbolic Trigonometric Functions
>>> linalg.sinhm(D)
Hyperbolic matrix sine
>>> linalg.coahm(D)
Hyperbolic matrix cosine
>>> linalg.tanhm(A)
Hyperbolic matrix tangent
Matrix Sign Function
>>> np.signm(A)
Matrix sign function
Matrix Square Root
>>> linalg.sqrtm(A)
Matrix square root
Arbitrary Functions
>>> linalg.funm(A, lambda x: x*x)
Evaluate matrix function
```

Basic Matrix Routines

```
Inverse
>>> A.I
Inverse
>>> linalg.inv(A)
Inverse
>>> A.I
Inverse
>>> A.H
Transpose matrix
>>> A.H
Conjugate transpose
Trace
>>> np.trace(A)
Trace
Norm
>>> linalg.norm(A)
Frobenius norm
>>> linalg.norm(A,1)
L1 norm (max column sum)
>>> linalg.norm(A,np.inf)
Linf norm (max row sum)
Rank
>>> np.linalg.matrix_rank(C)
Matrix rank
Determinant
>>> linalg.det(A)
Determinant
Solving linear problems
>>> linalg.solve(A,b)
Solver for dense matrices
>>> E = np.mat(A).T
>>> linalg.lstsq(F,E)
Least-squares solution to linear matrix equation
Generalized inverse
>>> linalg.pinv(C)
Compute the pseudo-inverse of a matrix (least-squares solver)
>>> linalg.pinv2(C)
Compute the pseudo-inverse of a matrix (SVD)
```

Creating Sparse Matrices

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

Sparse Matrix Routines

```
Inverse
>>> sparse.linalg.inv(I)
Inverse
Norm
>>> sparse.linalg.norm(I)
Norm
Solving linear problems
>>> sparse.linalg.spsolve(H,I)
Solver for sparse matrices
```

Sparse Matrix Functions

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

Asking For Help

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

Decompositions

```
Eigenvalues and Eigenvectors
>>> la, v = linalg.eig(A)
Solve ordinary or generalized eigenvalue problem for square matrix
>>> l1, l2 = la
Unpack eigenvalues
>>> v[:,0]
First eigenvector
>>> v[:,1]
Second eigenvector
>>> linalg.eigvals(A)
Unpack eigenvalues
Singular Value Decomposition
>>> U,s,Vh = linalg.svd(B)
Singular Value Decomposition (SVD)
>>> M,N = B.shape
>>> Sig = linalg.diagsvd(s,M,N)
Construct sigma matrix in SVD
LU Decomposition
>>> P,L,U = linalg.lu(C)
LU Decomposition
```

Sparse Matrix Decompositions

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

DataCamp

Learn Python for Data Science Interactively



Python For Data Science Cheat Sheet

Matplotlib

Learn Python Interactively at www.DataCamp.com



Matplotlib

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



1 Prepare The Data

Also see Lists & NumPy

1D Data

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

2D Data or Images

```
>>> data = 2 * np.random.random((10, 10))
>>> data2 = 3 * np.random.random((10, 10))
>>> Y, X = np.meshgrid([3:11*0.5], [-3:3*0.5])
>>> U = -1 - X**2 + Y
>>> V = 1 + X - Y**2
>>> from matplotlib.cbook import get_sample_data
>>> img = np.load(get_sample_data('axes_grid/bivariate_normal.npy'))
```

2 Create Plot

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

Figure

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

Axes

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

```
>>> fig.add_axes()
>>> ax1 = fig.add_subplot(221) # row-col-num
>>> ax3 = fig.add_subplot(212)
>>> fig3, axes = plt.subplots(ncrow=2, ncol=2)
>>> fig4, axes2 = plt.subplots(ncols=3)
```

3 Plotting Routines

1D Data

```
>>> fig, ax = plt.subplots()
>>> lines = ax.plot(x,y)
>>> ax.scatter(x,y)
>>> axes[0,0].bar([1,2,3],[3,4,5])
>>> axes[1,0].barh([0.5,1,2.5],[0,1,2])
>>> axes[1,1].axhline(0.45)
>>> axes[0,1].axvline(0.65)
>>> ax.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 vertical rectangles (constant width)
Plot horizontal rectangles (constant height)
Draw a horizontal line across axes
Draw a vertical line across axes
Draw filled polygons
Fill between y-values and 0

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)
```

Add an arrow to the axes
Plot a 2D field of arrows
Plot a 2D field of arrows

Data Distributions

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

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

2D Data or Images

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

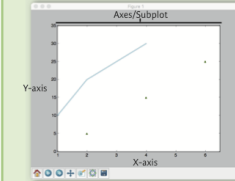
Colormapped or RGB arrays

```
>>> axes2[0].pcolor(data2)
>>> axes2[0].pcolormesh(data)
>>> CS = plt.contour(Y,X,U)
>>> axes2[2].contourf(data1)
>>> axes2[2] = ax.clabel(CS)
```

Pseudocolor plot of 2D array
Pseudocolor plot of 2D array
Plot contours
Plot filled contours
Label a contour plot

Plot Anatomy & Workflow

Plot Anatomy



Workflow

The basic steps to creating plots with matplotlib are:

- 1 Prepare data
- 2 Create plot
- 3 Plot
- 4 Customize plot
- 5 Save plot
- 6 Show plot

```
>>> import matplotlib.pyplot as plt
>>> x = [1,2,3,4]
>>> y = [10,20,25,30]
>>> fig = plt.figure()
>>> ax = fig.add_subplot(111)
>>> ax.plot(x, y, color='lightblue', linewidth=3)
>>> ax.scatter(x, y, color='darkgreen',
>>>           marker='^')
>>> ax.set_xlim(1, 6.5)
>>> plt.savefig('foo.png')
>>> plt.show()
```

4 Customize Plot

Colors, Color Bars & Color Maps

```
>>> plt.plot(x, x, x, x**2, x, x**3)
>>> ax.plot(x, y, alpha=0.4)
>>> ax.plot(x, y, c='r')
>>> fig.colorbar(im, orientation='horizontal')
>>> im = ax.imshow(img, cmap='seismic')
```

Markers

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

Linestyles

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

Text & Annotations

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

Mathtext

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

Limits, Legends & Layouts

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

Legends

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

Ticks

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

Subplot Spacing

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

Axis Spines

```
>>> ax1.spines['top'].set_visible(False)
>>> ax1.spines['bottom'].set_visible(True)
```

Add padding to a plot
Set the aspect ratio of the plot
Set limits for x and y-axis
Set limits for x-axis
Set a title and x and y-axis labels
No overlapping plot elements
Manually set x-ticks
Make y-ticks longer and go in and out
Adjust the spacing between subplots
Fit subplot(s) in to the figure area
Make the top axis line for a plot invisible
Make the bottom axis line outward

5 Save Plot

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

6 Show Plot

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

Close & Clear

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

Clear an axis
Clear the entire figure
Close a window

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, y = iris.data[0:,1:4], 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(['W', 'W', 'F', 'F', 'M', 'F', 'M', 'M', 'F', '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
>>> scaler = Normalizer().fit(X_train)
>>> normalized_X = scaler.transform(X_train)
>>> normalized_X_test = scaler.transform(X_test)
```

Binarization

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

Create Your Model

Supervised Learning Estimators

Linear Regression

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

Support Vector Machines (SVM)

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

Naive Bayes

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

KNN

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

Unsupervised Learning Estimators

Principal Component Analysis (PCA)

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

K Means

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

Model Fitting

Supervised Learning

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

Fit the model to the data

Unsupervised Learning

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

Fit the model to the data
Fit to data, then transform it

Prediction

Supervised Estimators

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

Predict labels
Predict labels
Estimate probability of a label

Unsupervised Estimators

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

Predict labels in clustering algos

Evaluate Your Model's Performance

Classification Metrics

Accuracy Score

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

Estimator score method
Metric scoring functions
Precision, recall, f1-score and support

Classification Report

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

Confusion Matrix

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

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
>>> metrics.mean_squared_error(y_true, y_pred)
>>> mean_squared_error(y_true, y_pred)
```

R² Score

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

Clustering Metrics

Adjusted Rand Index

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

Homogeneity

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

V-measure

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

Cross-Validation

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

Tune Your Model

Grid Search

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

Randomized Parameter Optimization

```
>>> from sklearn.grid_search import RandomizedSearchCV
>>> params = {'n_neighbors': range(1,5),
>>>          'weights': ["uniform", "distance"]}
>>> 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_)
```














DataCamp

Learn Python for Data Science Interactively

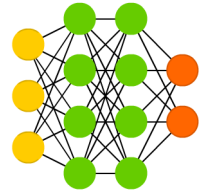


A mostly complete chart of Neural Networks

©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

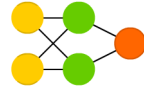
Deep Feed Forward (DFF)



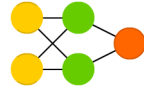
Perceptron (P)



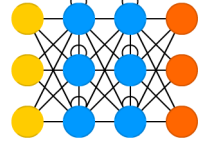
Feed Forward (FF)



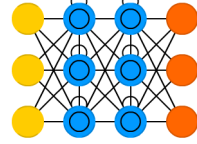
Radial Basis Network (RBF)



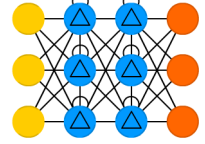
Recurrent Neural Network (RNN)



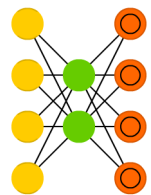
Long / Short Term Memory (LSTM)



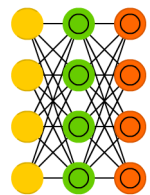
Gated Recurrent Unit (GRU)



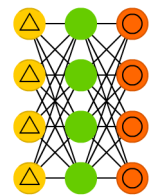
Auto Encoder (AE)



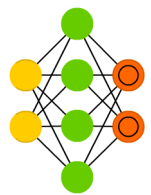
Variational AE (VAE)



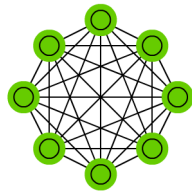
Denoising AE (DAE)



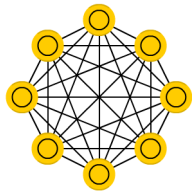
Sparse AE (SAE)



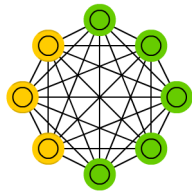
Markov Chain (MC)



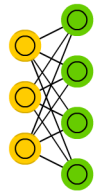
Hopfield Network (HN)



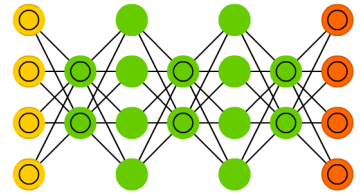
Boltzmann Machine (BM)



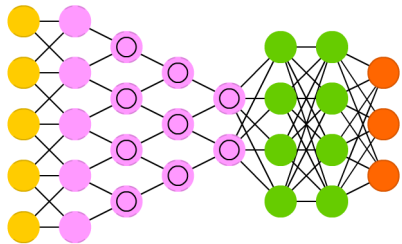
Restricted BM (RBM)



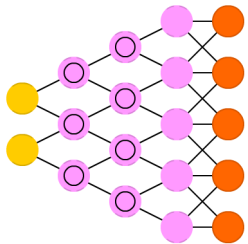
Deep Belief Network (DBN)



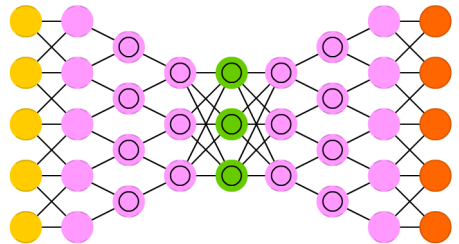
Deep Convolutional Network (DCN)



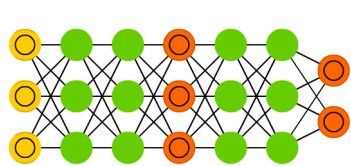
Deconvolutional Network (DN)



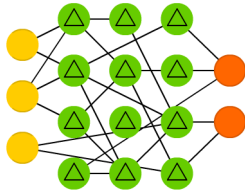
Deep Convolutional Inverse Graphics Network (DCIGN)



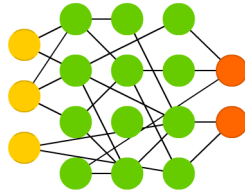
Generative Adversarial Network (GAN)



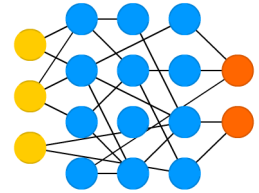
Liquid State Machine (LSM)



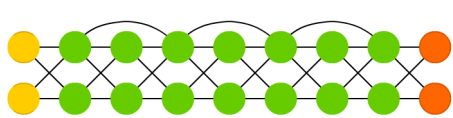
Extreme Learning Machine (ELM)



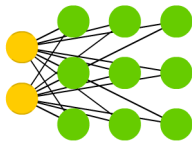
Echo State Network (ESN)



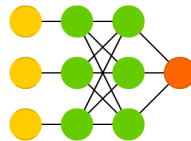
Deep Residual Network (DRN)



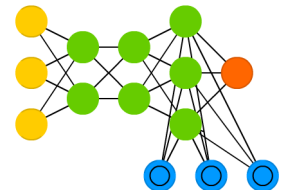
Kohonen Network (KN)



Support Vector Machine (SVM)



Neural Turing Machine (NTM)



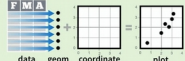
Data Visualization with ggplot2

Cheat Sheet

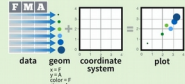


Basics

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 **qplot()** or **ggplot()**

qplot(x = cty, y = 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 = cty, y = hwy))
Begins a plot that you finish by adding layers to. No defaults, but provides more control than qplot().

ggplot(mpg, aes(hwy, cty)) + geom_point(aes(color = cyl)) + geom_smooth(method = "lm") + coord_cartesian() + scale_color_gradient() + theme_bw()
add layers, elements with +
layer = geom + default stat + layer specific mappings
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 for file extension.

RStudio® is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 844-448-1212 • rstudio.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
b + geom_area(aes(y = ..density..), stat = "bin")
- a + geom_density(kernel = "gaussian")**
x, y, alpha, color, fill, linetype, size, weight
b + geom_density(aes(y = ..county..))
- a + geom_dotplot()**
x, y, alpha, color, fill
- a + geom_freqpoly()**
x, y, alpha, color, linetype, size
b + geom_freqpoly(aes(y = ..density..))
- a + geom_histogram(binwidth = 5)**
x, y, alpha, color, fill, linetype, size, weight
b + geom_histogram(aes(y = ..density..))

Discrete

b <- ggplot(mpg, aes(fill))

- b + geom_bar()**
x, alpha, color, fill, linetype, size, weight

Graphical Primitives

c <- ggplot(mpg, 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(linetype = "butt", linejoin = "round", linemiter = 1)**
x, y, alpha, color, linetype, size
- d + geom_ribbon(aes(ymin = unemploy - 900, ymax = unemploy + 900))**
x, y, 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, y, alpha, color, fill, linetype, size
- e + geom_rect(aes(xmin = long, ymin = lat, xmax = long + delta_long, ymax = lat + delta_lat))**
x, y, alpha, color, fill, linetype, size

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
- f + geom_text(aes(label = cty))**
x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

Discrete X, Continuous Y

g <- ggplot(mpg, aes(class, hwy))

- g + geom_bar(stat = "identity")**
x, y, alpha, color, fill, linetype, size, weight
- g + geom_boxplot()**
lower, middle, upper, x, ymax, ymin, alpha, color, fill, linetype, shape, size, weight
- g + geom_dotplot(binaxis = "y", stackdir = "center")**
x, y, alpha, color, fill
- g + geom_violin(scale = "area")**
x, y, alpha, color, fill, linetype, size, weight

Discrete X, Discrete Y

h <- ggplot(diamonds, aes(cut, color))

- h + geom_jitter()**
x, y, alpha, color, fill, shape, size

Continuous Bivariate Distribution

i <- ggplot(movies, aes(date, rating))

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

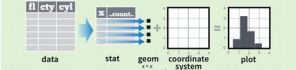
sealsz <- with(seals, sqrt(delta_long^2 + delta_lat^2))
m <- ggplot(seals, aes(long, lat))

- 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

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

Stats - An alternative way to build a layer

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")**

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

a + stat_bin(binwidth = 1, origin = 10) 1D distributions
x, y | count, ..ncount.., density, ..density..
a + stat_bin2d(binwidth = 1, bins = "xy")
x, y | count, ..ncount..
a + stat_density(adjust = 1, kernel = "gaussian")
x, y | count, ..density.., scaled..

f + stat_bin2d(bins = 30, drop = TRUE) 2D distributions
x, y, fill | count, ..density..
f + stat_binhex(bins = 30)
x, y, fill | count, ..density..
f + stat_density2d(contour = TRUE, n = 100)
x, y, color, size | ..level..

m + stat_contour(aes(z = z)) 3 Variables
x, y, z | order | ..level..
m + stat_spoke(aes(radius = z, angle = z))
angle, radius, x, yend, ..xend.., ..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_density2d(adjust = 1, kernel = "gaussian", scale = "area")
x, y | density, ..scaled..count.., ..n.., ..violinwidth.., width..

f + stat_cdf(n = 40) Functions
x, y | ..x.., ..y..
f + stat_quantile(quantiles = 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.., ..ymin.., ..ymax..

ggplot() + stat_function(aes(x = 3.3)) General Purpose
fun = dnorm, n = 101, args = list(sd = 0.5)
x | y..
f + stat_identity()
ggplot() + stat_qq(aes(sample = 1:100), distribution = qt, dparams = list(df = 5))
sample, x, y | ..x.., ..y..
f + stat_sum()
x, y, size | ..size..
f + stat_summary(fun.data = "mean_cl_boot")
f + stat_unique()

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 = fi))
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 legend/axis 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(values = c()) - map discrete values to manually chosen visual values

X and Y location scales

Use with x or y aesthetics (x shown here)

scale_x_date(labels = date_format("%m/%d"), 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_reverse() - Reverse direction of x axis
scale_x_sqrt() - Plot x on square root scale

Color and fill scales

n <- b + geom_bar(aes(fill = fi)) Discrete
n | scale_fill_brewer(palette = "Blues")
For palette choices: library(RColorBrewer) display(brewer.all)
n | scale_fill_grey(start = 0.2, end = 0.8, na.value = "red")
o <- a + geom_dotplot(aes(fill = ..x..)) Continuous
o | scale_fill_gradient(low = "red", high = "yellow")
o | scale_fill_gradient2(low = "red", high = "blue", mid = "white", midpoint = 25)
Also: rainbow(), heat.colors(), topo.colors(), cm.colors(), RColorBrewer:brewer.pall()

Shape scales

p <- f + geom_point(aes(shape = fi)) Manual shape values
p | scale_shape(solid = FALSE)
p | scale_shape_manual(values = c(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))
Shapes values shown in chart on right

Size scales

q <- f + geom_point(aes(size = ci)) Manual size values
q | scale_size_area(max = 5)
Value mapped to area of circle (not radius)

Coordinate Systems

- r <- b + geom_bar()**
- r + coord_cartesian(xlim = c(0, 5))**
xlim, ylim
The default cartesian coordinate system
- r + coord_fixed(ratio = 1/2)**
ratio, xlim, ylim
Cartesian coordinates with fixed aspect ratio between x and y units
- r + coord_flip()**
xlim, ylim
Flipped Cartesian coordinates
- r + coord_polar(theta = "x", direction = 1)**
theta, start, direction
Polar coordinates
- r + coord_trans(ytrans = "sqrt")**
xtrans, ytrans, xlim, ylim
Transformed cartesian coordinates. Set extras and strains to the name of a window function.

z + coord_map(projection = "ortho", orientation = c(41, -74, 0))
projection, orientation, xlim, ylim
Map projections from the mapproj package (mercator (default), azequialarea, lagrange, etc.)

Position Adjustments

- Position adjustments determine how to arrange geoms that would otherwise occupy the same space.
- s <- ggplot(mpg, aes(fi, fill = drv))**
- s + geom_bar(position = "dodge")**
Arrange elements side by side
- s + geom_bar(position = "fill")**
Stack elements on top of one another, normalize height
- s + geom_bar(position = "stack")**
Stack elements on top of one another
- f + geom_point(position = "jitter")**
Add random noise to X and Y position of each element to avoid overplotting

Each position adjustment can be recast as a function with manual **width** and **height** arguments

s + geom_bar(position = position_dodge(width = 1))

- f + theme_bw()** White background with grid lines
 - f + theme_classic()** White background no gridlines
 - f + theme_grey()** Grey background (default theme)
 - f + theme_minimal()** Minimal theme
- ggthemes - Package with additional ggplot2 themes

Faceting

Facets divide a plot into subplots based on the values of one or more discrete variables.

t <- ggplot(mpg, aes(cty, hwy)) + geom_point()

- t + facet_grid(. ~ fi)**
facet into columns based on fi
- t + facet_grid(year ~ .)**
facet into rows based on year
- t + facet_grid(year ~ fi)**
facet into both rows and columns
- t + facet_wrap(~ fi)**
wrap facets into a rectangular layout

Set **scales** to let axis limits vary across facets

- t + facet_grid(y ~ x, scales = "free")**
x and y axis limits adjust to individual facets
- "free_x"** - x axis limits adjust
- "free_y"** - y axis limits adjust

Set **labeller** to adjust facet labels

- t + facet_grid(. ~ fi, labeller = label_both)**
fi: c fi: d fi: e fi: p fi: r
- t + facet_grid(. ~ fi, labeller = label_bquote(alpha ~ .))**
α' α' α' α' α'
- t + facet_grid(. ~ fi, labeller = label_parsed)**
c d e p r

Labels

t + ggtitle("New Plot Title")
Add a main title above the plot
t + xlab("New X label")
Change the label on the X axis
t + ylab("New Y label")
Change the label on the Y axis
t + labs(title = "New title", x = "New x", y = "New y")
All of the above

Use scale functions to update legend labels

Legends

t + theme(legend.position = "bottom")
Place legend at "bottom", "top", "left", or "right"
t + guides(color = "none")
Set legend type for each aesthetic: color, legend, or none (no legend)
t + scale_fill_discrete(name = "Title", labels = c("A", "B", "C"))
Set legend title and labels with a scale function.

Zooming

Without clipping (preferred)
t + coord_cartesian(xlim = c(0, 100), ylim = c(10, 20))
With clipping (removes unseen data points)
t + xlim(0, 100) + ylim(10, 20)
t + scale_x_continuous(limits = c(0, 100))
t + scale_y_continuous(limits = c(0, 100))

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

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