

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)
>>> X_train3 = to_categorical(x_train3, num_classes)
>>> Y_test3 = to_categorical(y_test3, num_classes)
```

Model Architecture

Sequential Model

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

Multilayer Perceptron (MLP)

Binary Classification

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

Multi-Class Classification

```
>>> from keras.layers import Dropout
>>> model.add(Dense(512, activation='relu', input_shape=(784,)))
>>> model.add(Dropout(0.2))
>>> model.add(Dense(512, activation='relu'))
>>> model.add(Dropout(0.2))
>>> model.add(Dense(10, activation='softmax'))
```

Regression

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

Convolutional Neural Network (CNN)

```
>>> from keras.layers import Activation, Conv2D, MaxPooling2D, Flatten
>>> model2.add(Conv2D(32, (3, 3), padding='same', input_shape=x_train.shape[1:]))
>>> model2.add(Activation('relu'))
>>> model2.add(Conv2D(32, (3, 3)))
>>> model2.add(Activation('relu'))
>>> model2.add(MaxPooling2D(pool_size=(2, 2)))
>>> model2.add(Dropout(0.25))
>>> model2.add(Conv2D(64, (3, 3), padding='same'))
>>> model2.add(Activation('relu'))
>>> model2.add(Conv2D(64, (3, 3)))
>>> model2.add(Activation('relu'))
>>> model2.add(MaxPooling2D(pool_size=(2, 2)))
>>> model2.add(Dropout(0.25))
>>> model2.add(Flatten())
>>> model2.add(Dense(512))
>>> model2.add(Activation('relu'))
>>> model2.add(Dropout(0.5))
>>> model2.add(Dense(num_classes))
>>> model2.add(Activation('softmax'))
```

Recurrent Neural Network (RNN)

```
>>> from keras.layers import Embedding, LSTM
>>> model3.add(Embedding(20000, 128))
>>> model3.add(LSTM(128, dropout=0.2, recurrent_dropout=0.2))
>>> model3.add(Dense(1, activation='sigmoid'))
```

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

```
>>> model.compile(optimizer='adam',
>>>               loss='binary_crossentropy',
>>>               metrics=['accuracy'])
MLP: Binary Classification
>>> model.compile(optimizer='rmsprop',
>>>               loss='categorical_crossentropy',
>>>               metrics=['accuracy'])
MLP: Multi-Class Classification
MLP: Regression
>>> model.compile(optimizer='rmsprop',
>>>               loss='mse',
>>>               metrics=['mse'])
Recurrent Neural Network
>>> model.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=0.001, 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

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
>>> 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
Get one element
>>> df[1:]
>>> Country Capital Population
>>> 1 India New Delhi 1303171035
>>> 2 Brazil Brasilia 207847528
Get subset of a DataFrame
```

Selecting, Boolean Indexing & Setting

```
By Position
>>> df.iloc[0], (0)
>>> 'Belgium'
Select single value by row & column
>>> df.iat[0], (0)
>>> 'Belgium'
By Label
>>> df.loc[0], ['Country']
>>> 'Belgium'
Select single value by row & column labels
>>> df.at[0], ['Country']
>>> 'Belgium'
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[:, '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 > s where value is < or =
Use filter to adjust DataFrame
Setting
>>> s['a'] = 6
Set index a of Series to 6
```

Dropping

```
>>> s.drop(['a', 'c'])
>>> df.drop('Country', axis=1)
Drop values from rows (axis=0)
Drop values from columns (axis=1)
```

Sort & Rank

```
>>> df.sort_index(by='Country')
>>> s.sort()
>>> df.rank()
Sort by row or column index
Sort a series by its values
Assign ranks to entries
```

Retrieving Series/DataFrame Information

```
Basic Information
>>> df.shape
>>> df.index
>>> df.columns
>>> df.info()
>>> df.count()
(rows, columns)
Describe index
Describe DataFrame columns
Info on DataFrame
Number of non-NA values
Summary
>>> df.sum()
>>> df.cumsum()
>>> df.min()/df.max()
>>> df.idxmin()/df.idxmax()
>>> df.describe()
>>> df.mean()
>>> df.median()
Sum of values
Cumulative sum of values
Minimum/maximum values
Minimum/Maximum index value
Series statistics
Mean of values
Median of values
```

Applying Functions

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

Data Alignment

```
Internal Data Alignment
NA values are introduced in the indices that don't overlap:
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
>>> s + s3
>>> a 10.0
>>> b NaN
>>> c 5.0
>>> d 7.0
Arithmetic Operations with Fill Methods
You can also do the internal data alignment yourself with the help of the fill methods:
>>> s.add(s3, fill_value=0)
>>> a 10.0
>>> b -5.0
>>> c 5.0
>>> d 7.0
>>> s.sub(s3, fill_value=2)
>>> s.div(s3, fill_value=4)
>>> s.mul(s3, fill_value=3)
```

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

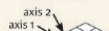
1D array



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))
```

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

I/O

Saving & Loading On Disk

```
>>> np.save('my_array', a)
>>> np savez('array.npy', 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.],
       [-3., -3., -3.]])
>>> np.subtract(a,b)
>>> b + a
array([[ 2.5,  4.,  6.],
       [ 5.,  7.,  9.]])
>>> np.add(b,a)
>>> a / b
array([[ 0.6666667,  1.,  1.],
       [ 0.25,  0.4,  0.5]])
>>> np.divide(a,b)
>>> a * b
array([[ 1.5,  4.,  9.],
       [ 4.,  10.,  18.]])
>>> np.multiply(a,b)
>>> np.exp(b)
>>> np.sqrt(b)
>>> np.sin(a)
>>> np.cos(b)
>>> np.log(a)
>>> a.dot(e)
array([[ 7.,  7.]])
```

Subtraction

Subtraction

Addition

Addition

Division

Division

Multiplication

Multiplication

Exponentiation

Square root

Print sines of an array

Element-wise cosine

Element-wise natural logarithm

Dot product

Comparison

```
>>> a == b
array([[False,  True,  True],
       [ True, 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.corrcoef()
>>> 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]
3
>>> b[1,2]
6.0
```

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

Slicing

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

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

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

Select all items at row 0 (equivalent to `b[0][:, :]`)
Same as `[1, 2, 3]`

Boolean Indexing

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

Reversed array `a`
Select elements from `a` less than 2

Fancy Indexing

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

Select elements from `a` less than 2
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, 4, 5)
>>> np.delete(a, [3])
```

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

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

Splitting Arrays

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

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

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], [4j,5j,6j]])
>>> c = np.array([[1,2,3,4], [5,6,7,8], [9,10,11,12]])
```

Index Tricks

```
>>> np.mgrid[0:5,0:5]          Create a dense meshgrid
>>> np.ogrid[0:2,0:2]        Create an open meshgrid
>>> np.r_[1,0]*5,-1:10]     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(c,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.cant['*'](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
>>> g[3:] += np.pi         (number of samples)
>>> np.unwrap(g)           Unwrap
>>> np.linspace(0,10,3)    Create an array of evenly spaced values (step size)
>>> 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 N-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

Creating Matrices

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

Basic Matrix Routines

Inverse	<code>linalg.inv(A)</code>	Inverse
Transpose	<code>A.T</code>	Transpose matrix Conjugate transposition
Trace	<code>np.trace(A)</code>	Trace
Norm	<code>linalg.norm(A)</code>	Frobenius norm L1 norm (max column sum) L inf norm (max row sum)
Rank	<code>np.linalg.matrix_rank(C)</code>	Matrix rank
Determinant	<code>linalg.det(A)</code>	Determinant
Solving linear problems	<code>linalg.solve(A,b)</code>	Solver for dense matrices Solver for sparse matrices Least-squares solution to linear matrix equation
Generalized inverse	<code>linalg.pinv(C)</code>	Compute the pseudo-inverse of a matrix (least-squares solver) Compute the pseudo-inverse of a matrix (SVD)

Matrix Functions

Addition	<code>np.add(A,D)</code>	Addition
Subtraction	<code>np.subtract(A,D)</code>	Subtraction
Division	<code>np.divide(A,D)</code>	Division
Multiplication	<code>np.multiply(D,A)</code>	Multiplication operator (Python 3) Dot product Vector dot product Inner product Outer product Tensor dot product Kronecker product
Exponential Functions	<code>linalg.exp(A)</code>	Matrix exponential Matrix exponential (Taylor Series) Matrix exponential (eigenvalue decomposition)
Logarithm Function	<code>linalg.log(A)</code>	Matrix logarithm
Trigonometric Functions	<code>linalg.sin(D)</code>	Matrix sine Matrix cosine Matrix tangent
Hyperbolic Trigonometric Functions	<code>linalg.sinh(D)</code>	Hyperbolic matrix sine Hyperbolic matrix cosine Hyperbolic matrix tangent
Matrix Sign Function	<code>np.sign(A)</code>	Matrix sign function
Matrix Square Root	<code>linalg.sqrtm(A)</code>	Matrix square root
Arbitrary Functions	<code>linalg.funm(A, lambda x: x*x)</code>	Evaluate matrix function

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[0 > 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.isapmatrix_csc(A) Identify sparse matrix
```

Sparse Matrix Routines

Inverse	<code>sparse.linalg.inv(I)</code>	Inverse
Norm	<code>sparse.linalg.norm(I)</code>	Norm
Solving linear problems	<code>sparse.linalg.solve(H,I)</code>	Solver for sparse matrices

Sparse Matrix Functions

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

Asking For Help

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

Decompositions

```
>>> la, v = linalg.eig(A)   Solve ordinary or generalized eigenvalue problem for square matrix
>>> U, S, Vt = la           Unpack eigenvalues
>>> U, S, Vh = la           First eigenvector
>>> linalg.eigvals(A)      Second eigenvector
>>> U, S, Vh = linalg.svd(B) Unpack eigenvalues
>>> M, N = B.shape          Singular Value Decomposition (SVD)
>>> Sig = linalg.diagsvd(S, N) Construct sigma matrix in SVD
>>> 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

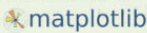
Python For Data Science Cheat Sheet

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

```
>>> import numpy as np
>>> x = np.linspace(0, 10, 100)
>>> y = np.cos(x)
>>> z = np.sin(x)
>>> data = 2 * np.random.random((10, 10))
>>> data2 = 3 * np.random.random((10, 10))
>>> Y, X = np.mgrid[0:10,0:10], -3:3:100j
>>> U = 1 - X**2 + Y**2
>>> V = 1 + X - Y**2
>>> from matplotlib.cbook import get_sample_data
>>> img = np.load(get_sample_data('mook_grid/bivariate_normal.npy'))
```

2 Create Plot

```
>>> import matplotlib.pyplot as plt
>>> fig = plt.figure()
>>> fig2 = plt.figure(figsize=plt.figaspect(2.0))
>>> 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_subplot(2,1,1)
>>> ax1 = fig.add_subplot(221) # row=col-num
>>> ax3 = fig.add_subplot(212)
>>> fig3, axes = plt.subplots(nrows=2,ncols=2)
>>> fig4, axes2 = plt.subplots(ncols=3)
```

3 Plotting Routines

```
>>> 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], [1,1,2])
>>> axes[1,1].axhline(0.45)
>>> axes[0,1].axvline(0.65)
>>> ax.fill_between(x,y,color='blue')
>>> ax.fill_between(x,y,color='yellow')
>>> fig, ax = plt.subplots()
>>> im = ax.imshow(img, cmap='gist_earth', interpolation='nearest', vmin=-2, vmax=2)
```

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([2,4,6], [15,25], color='darkgreen', marker='*')
>>> ax.set_xlim(1, 6.5)
>>> plt.savefig('foo.png')
>>> plt.show()
```

4 Customize Plot

```
>>> plt.plot(x, x, x**2, x, x**3)
>>> ax.plot(x, y, alpha = 0.4)
>>> ax.plot(x, y, c='k')
>>> fig.colorbar(im, orientation='horizontal')
>>> im = ax.imshow(img, cmap='seismic')
>>> fig, ax = plt.subplots()
>>> ax.scatter(x,y,marker="*")
>>> ax.plot(x,y,marker="o")
>>> 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.setp([lines,color='r',linewidth=4.0])
>>> ax.text(1, -2.1, "Example Graph", style='italic')
>>> ax.annotate("Sine", xy=(8, 0), xytext=(10, 5), textcoords='data', arrowprops=dict(arrowstyle="->", connectionstyle="arc3"),)
```

```
>>> plt.title(r'$\sigma_i=15$', fontsize=20)
>>> 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)
>>> ax.set(title='An Example Axes', ylabel='Y-Axis', xlabel='X-Axis')
>>> ax.legend(loc='best')
>>> ax.xaxis.set(ticks=range(1,5), ticklabels=[3,100,-12,'foo'])
>>> ax.tick_params(axis='y', direction='inout', length=10)
>>> fig3.subplots_adjust(wspace=0.5, hspace=0.3, left=0.125, right=0.9, top=0.9, bottom=0.1)
>>> fig.tight_layout()
>>> ax1.spines['top'].set_visible(False)
>>> ax1.spines['bottom'].set_position('outward',10)
```

5 Save Plot

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

6 Show Plot

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

Close & Clear

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

DataCamp

Learn Python for Data Science Interactively

Python For Data Science Cheat Sheet

PySpark Basics

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



Spark

PySpark is the Spark Python API that exposes the Spark programming model to Python



Initializing Spark

SparkContext

```
>>> from pyspark import SparkContext
>>> sc = SparkContext(master='local[2]')
```

Inspect SparkContext

```
>>> sc.version          Retrieve SparkContext version
>>> sc.pythonVer       Retrieve Python version
>>> sc.master           Master URL to connect to
>>> sc.sparkHome       Path where Spark is installed on worker nodes
>>> sc.sparkUser       Retrieve name of the Spark User running SparkContext
>>> sc.appName         Return application name
>>> sc.applicationId   Retrieve application ID
>>> sc.defaultParallelism Return default level of parallelism
>>> sc.defaultMinPartitions Default minimum number of partitions for RDDs
```

Configuration

```
>>> from pyspark import SparkConf, SparkContext
>>> conf = SparkConf()
>>> conf.setMaster("local")
>>> conf.setAppName("My app")
>>> conf.setSparkExecutorMemory("1g")
>>> sc = SparkContext(conf = conf)
```

Using The Shell

In the PySpark shell, a special interpreter-aware SparkContext is already created in the variable called `sc`.

```
$ ./bin/spark-shell --master local[2]
$ ./bin/pyspark --master local[4] --py-files code.py
```

Set which master the context connects to with the `--master` argument, and add Python `.zip`, `.egg` or `.py` files to the runtime path by passing a comma-separated list to `--py-files`.

Loading Data

Parallelized Collections

```
>>> rdd = sc.parallelize(['a', '7'], ['a', '2'], ['b', '2'])
>>> rdd2 = sc.parallelize(['a', '2'], ['d', '1'], ['b', '1'])
>>> rdd3 = sc.parallelize(range(100))
>>> rdd4 = sc.parallelize(["a", "b", "c", "d", "e"], ["a", "b", "c", "d", "e"])
```

External Data

Read either one text file from HDFS, a local file system or or any Hadoop-supported file system URI with `textFile()`, or read in a directory of text files with `wholeTextFiles()`.

```
>>> textFile = sc.textFile("my/directory/*.txt")
>>> textFile2 = sc.wholeTextFiles("my/directory/")
```

Retrieving RDD Information

Basic Information

```
>>> rdd.getNumPartitions()
>>> rdd.count()
>>> rdd.countByKey()
>>> rdd.countByValue()
>>> rdd.collect()
>>> rdd.collectAndMap(lambda x: x*2)
>>> rdd3.sum()
>>> sc.parallelize([]).isEmpty()
```

List the number of partitions
Count RDD instances
Count RDD instances by key
Count RDD instances by value
Return (key,value) pairs as a dictionary
Sum of RDD elements
Check whether RDD is empty

Summary

```
>>> rdd3.max()
>>> rdd3.min()
>>> rdd3.mean()
>>> rdd3.stdev()
>>> rdd3.variance()
>>> rdd3.histogram()
>>> rdd3.stats()
```

Maximum value of RDD elements
Minimum value of RDD elements
Mean value of RDD elements
Standard deviation of RDD elements
Compute variance of RDD elements
Compute histogram by bins
Summary statistics (count, mean, stdev, max & min)

Applying Functions

```
>>> rdd.map(lambda x: x*(x[1],x[0]))
>>> rdd5 = rdd.flatMap(lambda x: x*(x[1],x[0]))
>>> rdd5.collect()
>>> rdd4.flatMapValues(lambda x: x)
>>> rdd4.collect()
```

Apply a function to each RDD element
Apply a function to each RDD element and flatten the result
Apply a flatMap function to each (key,value) pair of RDD without changing the keys

Selecting Data

```
>>> rdd.collect()
>>> rdd.take(2)
>>> rdd.first()
>>> rdd.top(2)
>>> rdd3.sample(False, 0.15, 81).collect()
```

Return a list with all RDD elements
Take first 2 RDD elements
Take first RDD element
Take top 2 RDD elements
Return sampled subset of RDD3

Filtering

```
>>> rdd4.filter(lambda x: "a" in x)
>>> rdd5.distinct().collect()
>>> rdd.keys().collect()
```

Filter the RDD
Return distinct RDD values
Return (key,value) RDD's keys

Iterating

```
>>> def g(x): print(x)
>>> rdd.foreach(g)
>>> rdd.foreachPartition(g)
```

Apply a function to all RDD elements

Reshaping Data

Reducing

```
>>> rdd.reduceByKey(lambda x,y: x+y)
>>> rdd.reduce(lambda a, b: a + b)
>>> rdd3.groupBy(lambda x: x % 2)
>>> rdd.aggregate((0,0), seqOp, combOp)
>>> rdd3.aggregateByKey((0,0), seqOp, combOp)
>>> rdd3.fold(0, add)
>>> rdd3.foldByKey(0, add)
>>> rdd3.keyBy(lambda x: x*x)
>>> rdd3.cartesian(rdd2).collect()
```

Merge the RDD values for each key
Merge the RDD values
Return RDD of grouped values
Group RDD by key

Aggregating

```
>>> seqOp = (lambda x,y: (x[0]+y,x[1]+1))
>>> combOp = (lambda x,y: (x[0]+y[0],x[1]+y[1]))
>>> rdd3.aggregate((0,0), seqOp, combOp)
>>> rdd3.aggregateByKey((0,0), seqOp, combOp)
>>> rdd3.fold(0, add)
>>> rdd3.foldByKey(0, add)
>>> rdd3.keyBy(lambda x: x*x)
>>> rdd3.cartesian(rdd2).collect()
```

Aggregate RDD elements of each partition and then the results
Aggregate values of each RDD key
Aggregate the elements of each partition, and then the results
Merge the values for each key
Create tuples of RDD elements by applying a function

Mathematical Operations

```
>>> rdd.subtract(rdd2)
>>> rdd2.subtractByKey(rdd)
>>> rdd.cartesian(rdd2).collect()
```

Return each RDD value not contained in RDD2
Return each (key,value) pair of RDD2 with no matching key in RDD1
Return the Cartesian product of RDD1 and RDD2

Sort

```
>>> rdd2.sortBy(lambda x: x[1])
>>> rdd2.sortByKey()
>>> rdd2.sortBy(lambda x: x[1])
```

Sort RDD by given function
Sort (key, value) RDD by key

Partitioning

```
>>> rdd.repartition(4)
>>> rdd.coalesce(1)
```

New RDD with 4 partitions
Decrease the number of partitions in the RDD to 1

Saving

```
>>> rdd.saveAsTextFile("rdd.txt")
>>> rdd.saveAsHadoopFile("hdfs://namenodehost/parent/child",
                        org.apache.hadoop.mapred.TextOutputFormat)
```

Stopping SparkContext

```
>>> sc.stop()
```

Execution

```
$ ./bin/spark-submit examples/src/main/python/pi.py
```

DataCamp
Learn Python for Data Science Interactively

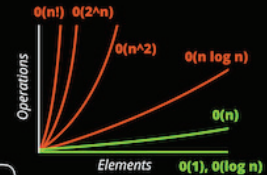
LEGEND

TIME Complexity VS. SPACE Complexity



<BIG-O-CHEATSHEET>

www.bigocheatsheet.com



DATA STRUCTURE	TIME Complexity				SPACE Complexity			
	Access	Search	Insertion	Deletion	Access	Search	Insertion	Deletion
Array	Good	Good	Good	Good	Good	Good	Good	Good
Stack	Good	Good	Good	Good	Good	Good	Good	Good
Queue	Good	Good	Good	Good	Good	Good	Good	Good
Singly-Linked List	Good	Good	Good	Good	Good	Good	Good	Good
Doubly-Linked List	Good	Good	Good	Good	Good	Good	Good	Good
Skip List	Good	Good	Good	Good	Good	Good	Good	Good
Hash Table	Good	Good	Good	Good	Good	Good	Good	Good
Binary Search Tree	Good	Good	Good	Good	Good	Good	Good	Good
Cartesian Tree	Good	Good	Good	Good	Good	Good	Good	Good
B-Tree	Good	Good	Good	Good	Good	Good	Good	Good
Red-Black Tree	Good	Good	Good	Good	Good	Good	Good	Good
Splay Tree	Good	Good	Good	Good	Good	Good	Good	Good
AVL Tree	Good	Good	Good	Good	Good	Good	Good	Good
KD Tree	Good	Good	Good	Good	Good	Good	Good	Good