

Python For Data Science Cheat Sheet

NumPy Basics

Learn Python for Data Science [Interactively](#) at [www.DataCamp.com](#)



NumPy

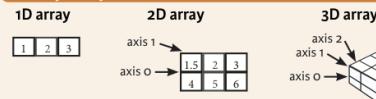
The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.

Use the following import convention:

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



NumPy Arrays



Creating Arrays

```
>>> a = np.array([1,2,3])
>>> b = np.array([(1,2,3), (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),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("myfile.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)
array([[ 0.66666667,  0.4,  0.5],
       [ 0.25,  0.5,  0.75]]])
>>> a / b
array([[ 1.5,  4.,  9.],
       [ 4., 10., 18.]])
>>> np.multiply(a,b)
>>> np.exp(b)
>>> np.sqrt(b)
>>> np.sin(a)
>>> np.cos(b)
>>> np.log(a)
>>> e.dot(f)
array([[ 7.,  7.],
       [ 7.,  7.]])
```

Subtraction

Addition

Division

Multiplication

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

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

Comparison

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

Element-wise comparison

Element-wise comparison

Array-wise comparison

Aggregate Functions

>>> a.sum() >>> a.min() >>> b.max(axis=0) >>> b.cumsum(axis=1) >>> a.mean() >>> b.median() >>> a.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]
1, 2, 3
4, 5, 6
```

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

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

Reversed array

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

Same as [1,:,:]

Boolean Indexing

```
>>> a[a<2]
array([1, 2, 3])
>>> b[1,0,1,0]
array([1, 0, 1, 0], [0, 1, 2, 0])
```

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

Array Manipulation

Transposing Array

```
>>> i = np.transpose(b)
>>> i.T
```

Permute array dimensions
Permute array dimensions

Changing Array Shape

```
>>> b.ravel()
```

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

Flatten the array

Reshape, but don't change data

Adding/Removing Elements

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

Return a new array with shape (2,6)

Append items to an array

Insert items in an array

Delete items from an array

Combining Arrays

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

Concatenate arrays

Stack arrays vertically (row-wise)

Stack arrays vertically (row-wise)

Stack arrays horizontally (column-wise)

Create stacked column-wise arrays

Create stacked column-wise arrays

Create stacked column-wise arrays

Splitting Arrays

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

Split the array horizontally at the 3rd index

Split the array vertically at the 2nd index

3. Pandas

DataCamp

Learn Python for Data Science [Interactively](#)



Data Wrangling

with pandas

Cheat Sheet

<http://pandas.pydata.org>

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.

	a	b	c
n	v		
d	1	4	7
e	2	5	10
f	3	6	11
g	4	7	12

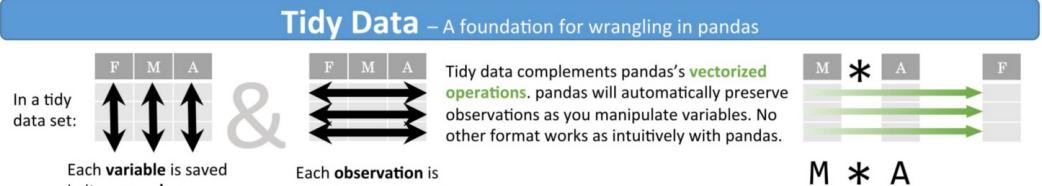
```
df = pd.DataFrame(
    {"a": [4, 5, 6],
     "b": [7, 8, 9],
     "c": [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [('d',1),('d',2),('e',2)],
        names=['n','v'])))
```

Create DataFrame with a MultiIndex

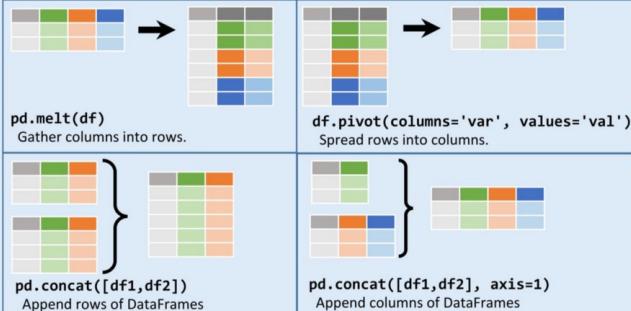
Method Chaining

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')
     )
```



Reshaping Data – Change the layout of a data set



```
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).

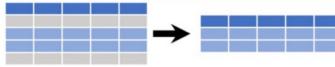
df.rename(columns = {'y':'year'})
Rename the columns of a DataFrame

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)



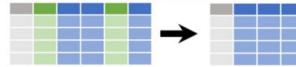
```
df[df.Length > 7]
Extract rows that meet logical criteria.

df.drop_duplicates()
Remove duplicate rows (only considers columns).

df.head(n)
Select first n rows.

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

Subset Variables (Columns)



```
df[['width', 'length', 'species']]
Select multiple columns with specific names.

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
'^^(?!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.
```

Logic in Python (and pandas)

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

<http://pandas.pydata.org/> This cheat sheet inspired by RStudio Data Wrangling Cheatsheet (<https://www.rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheatsheet.pdf>) Written by Irv Lustig, Princeton Consultants

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)
```

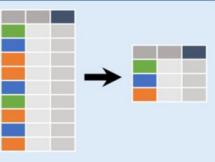


pandas provides a large set of **summary functions** that operate on different kinds of pandas objects (DataFrame columns, Series, GroupBy, Expanding and Rolling (see below)) and produce single values for each of the groups. When applied to a DataFrame, the result is returned as a pandas Series for each column. Examples:

```
sum()           Sum values of each object.
count()          Count non-NA/null values of each object.
median()         Median value of each object.
quantile([0.25,0.75]) Quantiles of each object.
apply(function)  Apply function to each object.
```

```
min()           Minimum value in each object.
max()           Maximum value in each object.
mean()          Mean value of each object.
var()            Variance of each object.
std()            Standard deviation of each object.
```

Group Data



```
df.groupby(by="col")
Return a GroupBy object, grouped by values in column named "col".
df.groupby(level="ind")
Return a GroupBy object, grouped by values in index level named "ind".
```

All of the summary functions listed above can be applied to a group. Additional GroupBy functions:

```
size()          Size of each group.
agg(function)   Aggregate group using function.
```

Windows

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

Handling Missing Data

```
df.dropna()
Drop rows with any column having NA/null data.
df.fillna(value)
Replace all NA/null data with value.
```

Make New Columns

```
df.assign(Area=lambda df: df.Length*df.Height)
Compute and append one or more new columns.
```

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

```
pd.qcut(df.col, n, labels=False)
Bin column into n buckets.
```



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

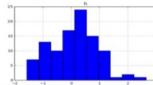
```
max(axis=1)      min(axis=1)
Element-wise max. Element-wise min.
clip(lower=-10,upper=10) abs()
Trim values at input thresholds Absolute value.
```

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

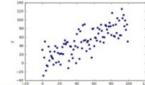
shift(1) Copy with values shifted by 1. rank(method='dense') Ranks with no gaps. rank(method='min') Ranks. Ties get min rank. rank(pct=True) Ranks rescaled to interval [0, 1]. rank(method='first') Ranks. Ties go to first value.	shift(-1) Copy with values lagged by 1. cumsum() Cumulative sum. cummax() Cumulative max. cummin() Cumulative min. cumprod() Cumulative product.
---	--

Plotting

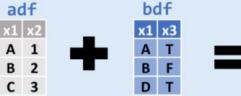
```
df.plot.hist()
Histogram for each column
```



```
df.plot.scatter(x='w',y='h')
Scatter chart using pairs of points
```



Combine Data Sets



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.
```



Set-like Operations

```
pd.merge(ydf, zdf)
Rows that appear in both ydf and zdf (Intersection).
```

```
pd.merge(ydf, zdf, how='outer')
Rows that appear in either or both ydf and zdf (Union).
```

```
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).
```

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

Index	A 3
	B -5
	C 7
	D 4

```
>>> 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

Index

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.iat[[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)]
>>> s[(s < -1) | (s > 2)]
>>> df[df['Population']>1200000000]
```

Series s where value is not > 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')
>>> s.order_()
>>> 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	(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*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([-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



4. Scipy

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([(1,5,2,3), (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.r_[3,[0]*5,-1:1:10j] 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.vstack((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 poly1d
>>> p = poly1d([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 obj to a data type
```

Other Useful Functions

```
>>> np.angle(b,deg=True) Return the angle of the complex argument
>>> np.linspace(0,np.pi,num=5) Create an array of evenly spaced values
... (number of samples)
>>> g [3:] += np.pi
>>> np.unwrap(g) Unwrap
>>> np.logspace(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
>>> Factorial Factorial
>>> misc.factorial(a) Combine N things taken at k time
>>> misc.comb(10,3,exact=True) Weights for N-point central derivative
>>> misc.central_diff_weights(3) Find the n-th derivative of a function at a point
>>> misc.derivative(myfunc,1.0)
```

Linear Algebra

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

```
>>> 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]])
```

Basic Matrix Routines

Inverse

```
>>> A.I
```

```
>>> linalg.inv(A)
```

Transpose

```
>>> A.T
```

```
>>> A.H
```

Trace

```
>>> np.trace(A)
```

Norm

```
>>> linalg.norm(A)
```

```
>>> linalg.norm(A,1)
```

```
>>> linalg.norm(A,np.inf)
```

Rank

```
>>> np.linalg.matrix_rank(C)
```

Determinant

```
>>> linalg.det(A)
```

Solving linear problems

```
>>> linalg.solve(A,b)
```

```
>>> E = np.mat(a).T
```

```
>>> linalg.lstsq(F,E)
```

Generalized inverse

```
>>> linalg.pinv(C)
```

```
>>> linalg.pinv2(C)
```

Frobenius norm

L1 norm (max column sum)

L inf norm (max row sum)

Matrix rank

Determinant

Solver for dense matrices

Solver for dense matrices

Least-squares solution to linear matrix equation

Compute the pseudo-inverse of a matrix

(least-squares solver)

Compute the pseudo-inverse of a matrix

(SVD)

Creating Sparse Matrices

Inverse

```
>>> F = np.eye(3, k=1)
```

```
>>> G = np.mat(np.identity(2))
```

```
>>> C[C > 0.5] = 0
```

```
>>> H = sparse.csr_matrix(C)
```

```
>>> I = sparse.csc_matrix(D)
```

```
>>> J = sparse.dok_matrix(A)
```

```
>>> E.todense()
```

```
>>> sparse.isspmatrix_csc(A)
```

Create a 2x2 identity matrix

Create a 2x2 identity matrix

Compressed Sparse Row matrix

Compressed Sparse Column matrix

Dictionary Of Keys matrix

Sparse matrix to full matrix

Identify sparse matrix

Sparse Matrix Routines

Inverse

```
>>> sparse.linalg.inv(I)
```

Norm

```
>>> sparse.linalg.norm(I)
```

Solving linear problems

```
>>> sparse.linalg.spsolve(H,I)
```

Inverse

Norm

Solver for sparse matrices

Sparse Matrix Functions

```
>>> sparse.linalg.expm(I)
```

Sparse matrix exponential

Eigenvalues and Eigenvectors

```
>>> la, v = linalg.eig(A)
```

Create ordinary or generalized eigenvalue problem for square matrix

Unpack eigenvalues

First eigenvector

Second eigenvector

Unpack eigenvalues

Singular Value Decomposition

```
>>> U,s,Vh = linalg.svd(B)
```

Singular Value Decomposition (SVD)

Construct sigma matrix in SVD

LU Decomposition

```
>>> P,L,U = linalg.lu(C)
```

LU Decomposition

Asking For Help

```
>>> help(scipy.linalg.diagsvd)
```

```
>>> np.info(np.matrix)
```

[Also see NumPy](#)

Matrix Functions

Addition

```
>>> np.add(A,D)
```

Subtraction

```
>>> np.subtract(A,D)
```

Division

```
>>> np.divide(A,D)
```

Multiplication

```
>>> A @ D
```

```
>>> np.multiply(D,A)
```

```
>>> np.dot(A,D)
```

```
>>> np.vdot(A,D)
```

```
>>> np.inner(A,D)
```

```
>>> np.outer(A,D)
```

```
>>> np.tensordot(A,D)
```

```
>>> np.kron(A,D)
```

Exponential Functions

```
>>> linalg.expm(A)
```

```
>>> np.expm2(A)
```

```
>>> np.expm3(D)
```

Logarithm Function

```
>>> linalg.logm(A)
```

Trigonometric Functions

```
>>> linalg.sinm(D)
```

```
>>> linalg.cosm(D)
```

```
>>> linalg.tanm(A)
```

Hyperbolic Trigonometric Functions

```
>>> linalg.sinhm(D)
```

```
>>> linalg.coshm(D)
```

```
>>> linalg.tanhm(A)
```

Matrix Sign Function

```
>>> np.signm(A)
```

Matrix Square Root

```
>>> linalg.sqrtm(A)
```

Arbitrary Functions

```
>>> linalg.funm(A, lambda x: x*x)
```

Decompositions

Eigenvalues and Eigenvectors

```
>>> la, v = linalg.eig(A)
```

Solve ordinary or generalized eigenvalue problem for square matrix

Unpack eigenvalues

First eigenvector

Second eigenvector

Unpack eigenvalues

Singular Value Decomposition

```
>>> U,s,Vh = linalg.svd(B)
```

Singular Value Decomposition

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](#)



5. Matplotlib

Python For Data Science Cheat Sheet

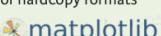
Matplotlib

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))  
>>> X, Y = np.mgrid[-3:3:100j, -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('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(nrows=2, ncols=2)  
>>> fig4, axes2 = plt.subplots(ncols=3)
```

3) Plotting Routines

1D Data

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

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

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

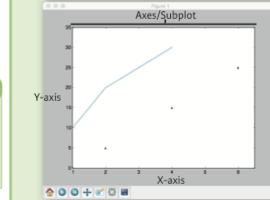
Data Distributions

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

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

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  
>>> 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 3,4  
>>> ax.scatter([2,4,6],  
[5,15,25],  
color='darkgreen',  
marker='^')  
>>> ax.set_xlim(1, 6.5) Step 4  
>>> plt.savefig('foo.png') Step 5  
>>> plt.show() Step 6
```

Figure

4) Customize Plot

Colors, Color Bars & Color Maps

```
>>> plt.plot(x, x, x**2, x, x**3)  
>>> ax.plot(x, y, alpha=.4)  
>>> ax.plot(x, c=x)  
>>> 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")  
>>> ax.setp(lines,color='r',linewidth=4.0)
```

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.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"))
```

Subplot Spacing

```
>>> fig3.subplots_adjust(wspace=0.5,
```

```
hspace=0.3, left=0.125,
```

```
right=0.9, top=0.9,
```

```
bottom=0.1)
```

Mathtext

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

Limits, Legends & Layouts

Limits & Autoscaling

```
>>> 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')
```

```
>>> ax.legend(loc='best')
```

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_position(('outward',10))
```

Add padding to a plot.
Set the aspect ratio of the plot to 1.
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.

Move the bottom axis line outward.

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()
```

Clear an axis.

Clear the entire figure.

Close a window.

6. Scikit-learn

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, y = iris.data[:, :2], iris.target
>>> X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=33)
>>> scaler = preprocessing.StandardScaler().fit(X_train)
>>> X_train = scaler.transform(X_train)
>>> X_test = scaler.transform(X_test)
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
>>> knn.fit(X_train, y_train)
>>> y_pred = knn.predict(X_test)
>>> accuracy_score(y_test, y_pred)
```

Loading The Data

Also see NumPy & Pandas

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

```
>>> import numpy as np
>>> X = np.random.random((10,5))
>>> y = np.array(['M', 'M', 'F', 'F', 'M', 'F', 'M', 'F', 'P', '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_train = scaler.transform(X_train)
>>> standardized_X_test = scaler.transform(X_test)
```

Normalization

```
>>> from sklearn.preprocessing import Normalizer
>>> scaler = Normalizer().fit(X_train)
>>> normalized_X_train = 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, y)
>>> 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

Encoding Categorical Features

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

Imputing Missing Values

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

Generating Polynomial Features

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

Evaluate Your Model's Performance

Classification Metrics

Accuracy Score
>>> from sklearn.metrics import accuracy_score
>>> from sklearn.metrics import accuracy_score
>>> accuracy_score(y_test, y_pred)

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))

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

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
>>> 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=5,
...                               n_iter=8,
...                               random_state=5)
>>> rsearch.fit(X_train, y_train)
>>> print(rsearch.best_score_)
```

DataCamp
Learn Python for Data Science interactively



7. Neural Networks Zoo

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

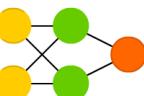
Perceptron (P)



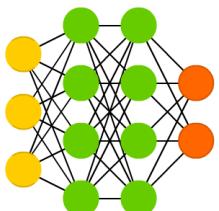
Feed Forward (FF)



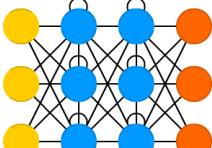
Radial Basis Network (RBF)



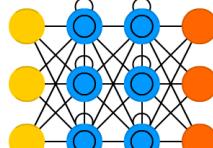
Deep Feed Forward (DFF)



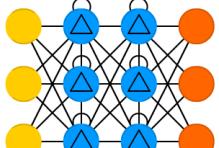
Recurrent Neural Network (RNN)



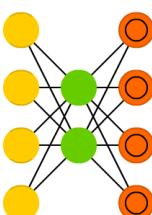
Long / Short Term Memory (LSTM)



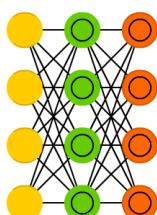
Gated Recurrent Unit (GRU)



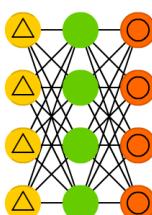
Auto Encoder (AE)



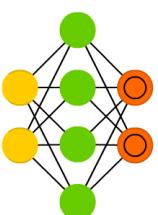
Variational AE (VAE)



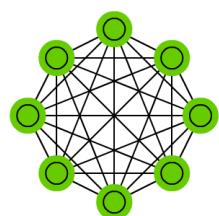
Denoising AE (DAE)



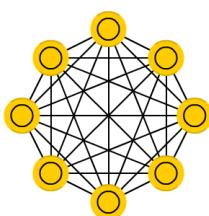
Sparse AE (SAE)



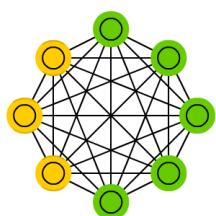
Markov Chain (MC)



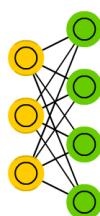
Hopfield Network (HN)



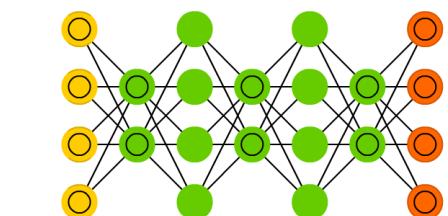
Boltzmann Machine (BM)



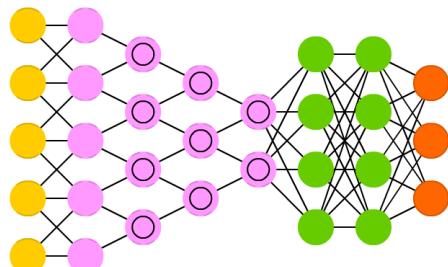
Restricted BM (RBM)



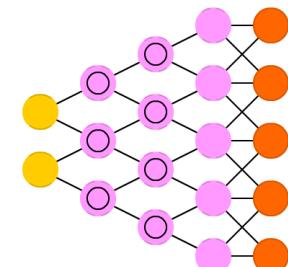
Deep Belief Network (DBN)



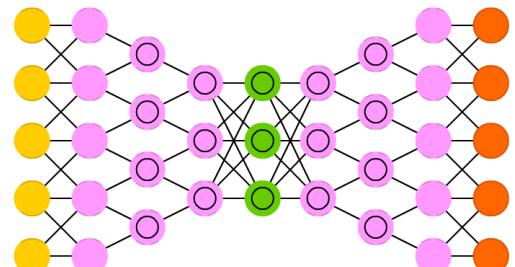
Deep Convolutional Network (DCN)



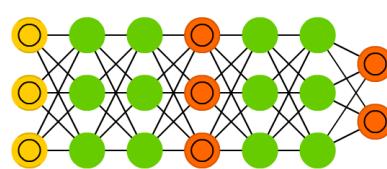
Deconvolutional Network (DN)



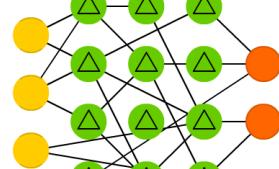
Deep Convolutional Inverse Graphics Network (DCIGN)



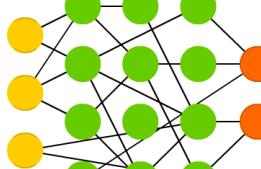
Generative Adversarial Network (GAN)



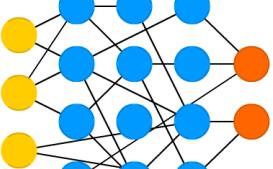
Liquid State Machine (LSM)



Extreme Learning Machine (ELM)



Echo State Network (ESN)



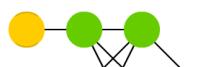
Deep Residual Network (DRN)



Kohonen Network (KN)

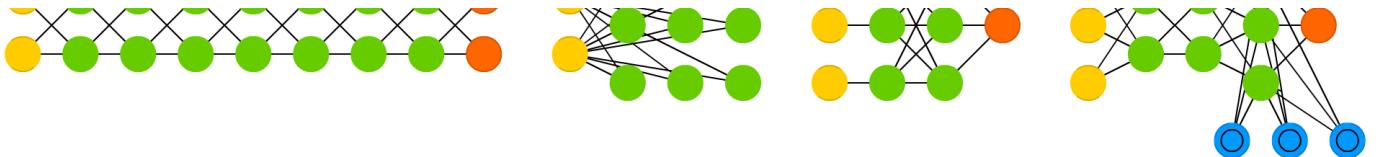


Support Vector Machine (SVM)



Neural Turing Machine (NTM)





8. ggplot2

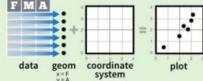
Data Visualization with ggplot2

Cheat Sheet

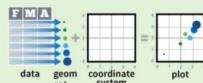


Basics

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 = cyl, 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 a new layer to a plot with a **geom_***() 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.

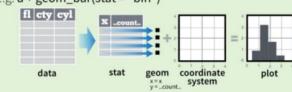
RStudio® is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 844-448-1212 • rstudio.com

Geoms - Use a geom to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.			
One Variable <ul style="list-style-type: none"> Continuous <ul style="list-style-type: none"> a + geom_area(stat = "bin") x, y, alpha, color, fill, linetype, size b + geom_area(aes(y = density..), stat = "bin") x, y, alpha, color, fill, linetype, size, weight b + geom_density(kernel = "gaussian") x, y, alpha, color, fill, linetype, size, weight b + geom_freqpoly(aes(y = ..density..)) x, y, alpha, color, linetype, size b + geom_histogram(binwidth = 5) x, y, alpha, color, fill, linetype, size, weight b + geom_histogram(aes(y = ..density..)) Discrete <ul style="list-style-type: none"> b <- ggplot(mpg, aes(f1)) b + geom_bar() x, alpha, color, fill, linetype, size, weight 		Two Variables <ul style="list-style-type: none"> Continuous X, Continuous Y <ul style="list-style-type: none"> 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 = cyl)) x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust 	
Graphical Primitives <ul style="list-style-type: none"> 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, 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, xend, y, yend, alpha, color, linetype, size e + geom_rect(aes(xmin = long, ymin = lat, xmax = long + delta_long, ymax = lat + delta_lat)) xmin, xmax, ymin, ymax, alpha, color, fill, linetype, size 		Continuous Bivariate Distribution <ul style="list-style-type: none"> i <- ggplot(movies, aes(year, rating)) i + geom_hex() x, y, alpha, colour, fill size 	
		Continuous Function <ul style="list-style-type: none"> 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 <ul style="list-style-type: none"> 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, y, ymax, ymin, alpha, color, linetype, size, width (also geom_errorbarh()) k + geom_linerange() x, y, ymin, ymax, alpha, color, linetype, size k + geom_pointrange() x, y, ymin, ymax, alpha, color, fill, linetype, shape, size 	
		Discrete X, Continuous Y <ul style="list-style-type: none"> 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, y, max, y, min, 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 <ul style="list-style-type: none"> h <- ggplot(diamonds, aes(cut, color)) h + geom_jitter() x, y, alpha, color, fill, shape, size 	
		Three Variables <ul style="list-style-type: none"> seals\$z <- 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_contour(aes(z = z)) x, y, z, alpha, colour, linetype, size, weight 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

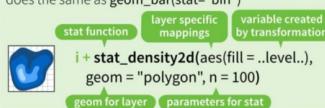
Stats - An alternative way to build a layer

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")`



```
a + stat_bin(binwidth = 1, origin = 10)          1D distributions
x, y, fill | .count..., density..., ndensity...
a + stat_bindot(binwidth = 1, binaxis = "x")      x, y, | .count..., .ncount...
a + stat_density(n = 100)                         x, y, | .count..., .scaled...
a + stat_density2d(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_hex(bins = 30)                          x, y, | .count..., .density...
f + stat_contour(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, xend, 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_density(adjust = 1, kernel = "gaussian", scale = "area")
x, y | .density..., .count..., .n..., .violinwidth..., .width...
```

```
f + stat_ecdf(n = 40)                           Functions
x, y | .x..., .y...
f + stat_quantile(qu quantiles = c(0.25, 0.5, 0.75), formula = y ~ log(x),
method = "rq")
x, y | .quantile..., .y...
f + stat_smooth(method = "auto", formula = y ~ x, se = TRUE, n = 80,
fullrange = FALSE, level = 0.95)
x, y | .se..., .x..., .y..., .ymin..., .ymax...
```

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

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

Scales

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



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/%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 x on log10 scale

`scale_x_reverse()` - Reverse direction of x axis

`scale_x_sqrt()` - Plot x on square root scale

Color and fill scales

Discrete

`n <- a + geom_bar(aes(fill = f))`

`n + scale_fill_brewer(palette = "Blues")`
For palette choices: blues(RdBrewer::brewer.all)

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

Continuous

`o <- a + geom_dotplot(aes(fill = ...))`

`o + scale_fill_gradient(low = "red", high = "blue")`
For gradient choices: RdBu(RColorBrewer::RdBu)

`o + scale_fill_gradient2(low = "red", high = "blue", mid = "white", midpoint = 25)`
For gradient2 choices: RdGy(RColorBrewer::RdGy)

`o + scale_fill_gradientn(colors = terrain.colors(6))`
Alternative: heat colors, topo.colors(), cm.colors(), RColorBrewer::brewer.all()

Shape scales

Manual shape values

0	6	12	18	24
1	7	13	19	25
2	8	14	20	*
3	9	15	21	...
4	10	16	22	O
5	11	17	23	◇

Size scales

`q <- f + geom_point(aes(size = cyl))`

`q + scale_size_area(max = 6, value mapped to area of circle (not radius))`

Coordinate Systems

`r <- b + geom_bar()`

`r + coord_cartesian(xlim = c(0, 5), ylim = c(0, 5))`

The default cartesian coordinate system

`r + coord_fixed(ratio = 1/2)`

Cartesian coordinates with fixed aspect ratio between x and y units

`r + coord_flip()`

Flipped Cartesian coordinates

`r + coord_polar(theta = "x", direction = 1)`

Polar coordinates

`r + coord_trans(ytrans = "sqrt")`

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), azequalarea, lagrange, etc.)

Position Adjustments

Position adjustments determine how to arrange geoms that would otherwise occupy the same space.

`s <- ggplot(mpg, aes(f, 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))`

Themes

`r + theme_bw()`

White background with grid lines

`r + theme_classic()`

White background no gridlines

`r + theme_grey()`

Grey background (default theme)

`r + 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(. ~ fl)`

facet into columns based on fl

`t + facet_grid(year ~ .)`

facet into rows based on year

`t + facet_grid(year ~ fl)`

facet into both rows and columns

`t + facet_wrap(~ fl)`

wrap facets into a rectangular layout

Set `scales` to let axis limits vary across facets

`t + facet_grid(~ 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 `labeler` to adjust facet labels

`t + facet_grid(~ fl, labeler = label_both)`

t: c	t: d	t: e	t: p	t: r
------	------	------	------	------

`t + facet_grid(~ fl, labeler = label_bquote(alpha ^ .(x)))`

a ^c	a ^d	a ^e	a ^p	a ^r
----------------	----------------	----------------	----------------	----------------

`t + facet_grid(~ fl, labeler = 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: colorbar, 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)) + scale_y_continuous(limits = c(0, 100))`

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

