

SciPy 基础

金 林

中南财经政法大学统计系

jinlin82@qq.com

2020 年 2 月 1 日



- 1 INTRODUCTION
- 2 SUB-PACKAGES
- 3 SCIPY.STATS
- 4 SCIPY.LINALG
- 5 SCIPY.CLUSTER
- 6 SCIPY.SPATIAL
- 7 DISTANCE COMPUTATIONS (SCIPY.SPATIAL.DISTANCE)
- 8 MULTI-DIMENSIONAL IMAGE PROCESSING (SCIPY.NDIMIMAGE)



Facts

- ❶ Initial release: Around 2001
- ❷ Stable release: 0.18.1 / 22 September 2016
- ❸ Website: <http://www.scipy.org>
- ❹ History: http://scipy.github.io/old-wiki/pages/History_of_SciPy



What is SciPy?

- 1 SciPy is an open source Python library used for scientific computing and technical computing.
- 2 SciPy builds on the NumPy array object
- 3 and is part of the NumPy stack which includes tools like Matplotlib, pandas and SymPy.
- 4 There is an expanding set of scientific computing libraries that are being added to the NumPy stack every day.
- 5 The NumPy stack is also sometimes referred to as the SciPy stack.



What can SciPy do?

- ❶ It adds significant power to the interactive Python session by providing the user with high-level commands and classes for manipulating and visualizing data.
- ❷ With SciPy an interactive Python session becomes a data-processing and system-prototyping environment rivaling systems such as MATLAB, IDL, Octave, R-Lab, and SciLab.
- ❸ The additional benefit of basing SciPy on Python is that this also makes a powerful programming language available for use in developing sophisticated programs and specialized applications.
- ❹ Scientific applications using SciPy benefit from the development of additional modules in numerous niches of the software landscape by developers across the world.
- ❺ Everything from parallel programming to web and data-base subroutines and classes have been made available to the Python programmer.



- 1 INTRODUCTION
- 2 SUB-PACKAGES
- 3 SCIPY.STATS
- 4 SCIPY.LINALG
- 5 SCIPY.CLUSTER
- 6 SCIPY.SPATIAL
- 7 DISTANCE COMPUTATIONS (SCIPY.SPATIAL.DISTANCE)
- 8 MULTI-DIMENSIONAL IMAGE PROCESSING (SCIPY.NDIMIMAGE)



Sub-packages

- ❶ constants: physical constants and conversion factors
- ❷ cluster: hierarchical clustering, vector quantization, K-means
- ❸ fftpack: Discrete Fourier Transform algorithms
- ❹ integrate: numerical integration routines
- ❺ interpolate: interpolation tools
- ❻ io: data input and output
- ❼ lib: Python wrappers to external libraries
- ❽ linalg: linear algebra routines
- ❾ misc: miscellaneous utilities (e.g. image reading/writing)



Sub-packages

- ⑩ ndimage: various functions for multi-dimensional image processing
- ⑪ optimize: optimization algorithms including linear programming
- ⑫ signal: signal processing tools
- ⑬ sparse: sparse matrix and related algorithms
- ⑭ spatial: KD-trees, nearest neighbors, distance functions
- ⑮ special: special functions
- ⑯ stats: statistical functions
- ⑰ weave: tool for writing C/C++ code as Python multiline strings



Sub-package

- 1 Scipy sub-packages need to be imported separately, for example:

```
from scipy import stats
```



- 1 INTRODUCTION
- 2 SUB-PACKAGES
- 3 **SCIPY.STATS**
- 4 SCIPY.LINALG
- 5 SCIPY.CLUSTER
- 6 SCIPY.SPATIAL
- 7 DISTANCE COMPUTATIONS (SCIPY.SPATIAL.DISTANCE)
- 8 MULTI-DIMENSIONAL IMAGE PROCESSING (SCIPY.NDIMIMAGE)



Intro

- ① This module contains a large number of probability distributions as well as a growing library of statistical functions.
- ② Each univariate distribution is an instance of a subclass of `rv_continuous` (`rv_discrete` for discrete distributions):
 - ① `rv_continuous([momtype, a, b, xtol, ...])` A generic continuous random variable class meant for subclassing.
 - ② `rv_discrete([a, b, name, badvalue, ...])` A generic discrete random variable class meant for subclassing.
- ③ Over 80 continuous random variables (RVs) and 10 discrete random variables have been implemented using these classes.
- ④ Besides this, new routines and distributions can easily added by the end user.



常见连续分布

beta	A beta continuous random variable.
cauchy	A Cauchy continuous random variable.
chi2	A chi-squared continuous random variable.
expon	An exponential continuous random variable.
f	An F continuous random variable.
gamma	A gamma continuous random variable.
ncf	A non-central F distribution continuous random variable
nct	A non-central Student' s T continuous random variable
norm	A normal continuous random variable.
norminvgauss	A Normal Inverse Gaussian continuous random variable
pareto	A Pareto continuous random variable.
t	A Student' s T continuous random variable.
uniform	A uniform continuous random variable.



常见离散分布

bernoulli	A Bernoulli discrete random variable.
binom	A binomial discrete random variable.
geom	A geometric discrete random variable.
hypergeom	A hypergeometric discrete random variable.
nbinom	A negative binomial discrete random variable.
poisson	A Poisson discrete random variable.
randint	A uniform discrete random variable.



常见多元分布

<code>multivariate_normal</code>	A multivariate normal random variable.
<code>matrix_normal</code>	A matrix normal random variable.
<code>dirichlet</code>	A Dirichlet random variable.
<code>wishart</code>	A Wishart random variable.
<code>invwishart</code>	An inverse Wishart random variable.
<code>multinomial</code>	A multinomial random variable.



例子：正态分布

```

1 from scipy import stats
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 mean, var, skew, kurt=stats.norm.stats(moments="mvsk")
6 x=np.linspace(-3, 3, 100)
7 plt.plot(x, stats.norm.pdf(x), label='norm pdf')
8 plt.plot(x, stats.norm.pdf(x, 3, 2), label='norm pdf')
9
10 ### Freeze the distribution and display the frozen pdf
11 rv=stats.norm(3,2)
12 rv.ppf(0.5)
13 rv.pdf(3)
14 rv.cdf(10)
15
16 ### Generate random numbers
17 r = stats.norm.rvs(size=1000)
18 plt.hist(r, density=True, histtype='stepfilled', alpha=0.2)

```



例子：F 分布

```
1 from scipy import stats
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 mean, var, skew, kurt=stats.f.stats(3, 5, moments="mvsk")
6 x=np.linspace(0.01, 6, 100)
7 plt.plot(x, stats.f.pdf(x, 3, 5), label='norm pdf')
8 plt.plot(x, stats.f.pdf(x, 3, 2), label='norm pdf')
9
10 ### Freeze the distribution and display the frozen pdf
11 rv=stats.f(3,2)
12 rv.ppf(0.5)
13 rv.pdf(3)
14 rv.cdf(10)
15
16 ### Generate random numbers
17 r = stats.f.rvs(3,2,size=1000)
18 plt.hist(r, density=True, histtype='stepfilled', alpha=0.2)
```



例子：二项分布

```
1 from scipy import stats
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 fig, ax = plt.subplots(1, 1)
6
7 mean, var, skew, kurt = stats.binom.stats(5, 0.4, moments='mvsk
8     ')
9
10 x = np.arange(stats.binom.ppf(0.01, 5, 0.4), stats.binom.ppf
11     (0.99, 5, 0.4))
12
13 ax.plot(x, stats.binom.pmf(x, 5, 0.4), 'bo', ms=8, label='binom
14     pmf')
15
16 ax.vlines(x, 0, stats.binom.pmf(x, 5, 0.4), colors='b', lw=5,
17     alpha=0.5)
18
19 rv = stats.binom(5, 0.4)
20 ax.vlines(x, 0, rv.pmf(x), colors='k', linestyle='-', lw=1,
21     label='frozen pmf')
22
23 ax.legend(loc='best', frameon=False)
```



例子：多元正态分布

```
1 from scipy import stats
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 x = np.linspace(0, 5, 10, endpoint=False)
6 y = stats.multivariate_normal.pdf(x, mean=2.5, cov=0.5); y
7
8 fig1 = plt.figure()
9 ax = fig1.add_subplot(111)
10 ax.plot(x, y)
11
12 x, y = np.mgrid[-1:1:.01, -1:1:.01]
13 pos = np.dstack((x, y))
14 rv = stats.multivariate_normal([0.5, -0.2], [[2.0, 0.3], [0.3,
15         0.5]])
16 fig2 = plt.figure()
17 ax2 = fig2.add_subplot(111)
18 ax2.contourf(x, y, rv.pdf(pos))
```



Common Methods

<code>rvs(loc=0, scale=1, size=1)</code>	Random variates.
<code>pdf(x, loc=0, scale=1)</code>	Probability density function.
<code>logpdf(x, loc=0, scale=1)</code>	Log of the probability density function.
<code>cdf(x, loc=0, scale=1)</code>	Cumulative distribution function.
<code>logcdf(x, loc=0, scale=1)</code>	Log of the cumulative distributionfunction.
<code>sf(x, loc=0, scale=1)</code>	Survival function (also defined as 1 - cdf, but sf is some
<code>logsf(x, loc=0, scale=1)</code>	Log of the survival function.
<code>ppf(q, loc=0, scale=1)</code>	Percent point function (inverse of cdf —percentiles).
<code>isf(q, loc=0, scale=1)</code>	Inverse survival function (inverse of sf).
<code>moment(n, loc=0, scale=1)</code>	Non-central moment of order n
<code>stats(loc=0, scale=1, moments='mv')</code>	Mean('m'), variance('v'), skew('s'), and/or ku
<code>entropy(loc=0, scale=1)</code>	(Differential) entropy of the RV.
<code>fit(data, loc=0, scale=1)</code>	Parameter estimates for generic data.
<code>expect(func, args=(), **kwds)</code>	Expected value of a function (of one argument) with res
<code>median(loc=0, scale=1)</code>	Median of the distribution.
<code>mean(loc=0, scale=1)</code>	Mean of the distribution.
<code>var(loc=0, scale=1)</code>	Variance of the distribution.
<code>std(loc=0, scale=1)</code>	Standard deviation of the distribution.
<code>interval(alpha, loc=0, scale=1)</code>	Endpoints of the range that contains alpha percent of t



几个统计函数

`describe(a[, axis, ddof, bias])`
`gmean(a[, axis, dtype])`
`hmean(a[, axis, dtype])`
`kurtosis(a[, axis, fisher, bias])`
`kurtosistest(a[, axis])`
`mode(a[, axis, nan_policy])`
`moment(a[, moment, axis])`
`normaltest(a[, axis, nan_policy])`
`skew(a[, axis, bias, nan_policy])`
`skewtest(a[, axis, nan_policy])`
`kstat(data[, n])`
`kstatvar(data[, n])`
`variation(a[, axis, nan_policy])`
`find_repeats(arr)`

Compute several descriptive statistics of the passed array.
 Compute the geometric mean along the specified axis.
 Calculate the harmonic mean along the specified axis.
 Compute the kurtosis (Fisher or Pearson) of a dataset.
 Test whether a dataset has normal kurtosis.
 Return an array of the modal (most common) value in the pa
 Calculate the nth moment about the mean for a sample.
 Test whether a sample differs from a normal distribution.
 Compute the skewness of a data set.
 Test whether the skew is different from the normal distributio
 Return the nth k-statistic ($1 \leq n \leq 4$ so far).
 Returns an unbiased estimator of the variance of the k-statisti
 Compute the coefficient of variation, the ratio of the biased st
 Find repeats and repeat counts.



Fitting Distributions



Building Specific Distributions



Analysing One Sample



Comparing two samples



Kernel Density Estimation

`scipy.stats.gaussian_kde`

- 1 Kernel density estimation is a way to estimate the probability density function (PDF) of a random variable in a non-parametric way.
- 2 `gaussian_kde` works for both uni-variate and multi-variate data.
- 3 It includes automatic bandwidth determination.
- 4 The estimation works best for a unimodal distribution; bimodal or multi-modal distributions tend to be oversmoothed.



- 1 INTRODUCTION
- 2 SUB-PACKAGES
- 3 SCIPY.STATS
- 4 SCIPY.LINALG
- 5 SCIPY.CLUSTER
- 6 SCIPY.SPATIAL
- 7 DISTANCE COMPUTATIONS (SCIPY.SPATIAL.DISTANCE)
- 8 MULTI-DIMENSIONAL IMAGE PROCESSING (SCIPY.NDIMIMAGE)



Linear algebra functions

- 1 `numpy.linalg` for more linear algebra functions.
- 2 Note that although `scipy.linalg` imports most of them, identically named functions from `scipy.linalg` may offer more or slightly differing functionality.



- 1 INTRODUCTION
- 2 SUB-PACKAGES
- 3 SCIPY.STATS
- 4 SCIPY.LINALG
- 5 SCIPY.CLUSTER
- 6 SCIPY.SPATIAL
- 7 DISTANCE COMPUTATIONS (SCIPY.SPATIAL.DISTANCE)
- 8 MULTI-DIMENSIONAL IMAGE PROCESSING (SCIPY.NDIMIMAGE)



- 1 INTRODUCTION
- 2 SUB-PACKAGES
- 3 SCIPY.STATS
- 4 SCIPY.LINALG
- 5 SCIPY.CLUSTER
- 6 SCIPY.SPATIAL
- 7 DISTANCE COMPUTATIONS (SCIPY.SPATIAL.DISTANCE)
- 8 MULTI-DIMENSIONAL IMAGE PROCESSING (SCIPY.NDIMENSIONAL_IMAGE_PROCESSING)



- 1 INTRODUCTION
- 2 SUB-PACKAGES
- 3 SCIPY.STATS
- 4 SCIPY.LINALG
- 5 SCIPY.CLUSTER
- 6 SCIPY.SPATIAL
- 7 DISTANCE COMPUTATIONS (SCIPY.SPATIAL.DISTANCE)
- 8 MULTI-DIMENSIONAL IMAGE PROCESSING (SCIPY.NDIM)



- 1 INTRODUCTION
- 2 SUB-PACKAGES
- 3 SCIPY.STATS
- 4 SCIPY.LINALG
- 5 SCIPY.CLUSTER
- 6 SCIPY.SPATIAL
- 7 DISTANCE COMPUTATIONS (SCIPY.SPATIAL.DISTANCE)
- 8 MULTI-DIMENSIONAL IMAGE PROCESSING (SCIPY.NDIMIMAGE)

