

1 Probability and Statistics Functions

1.1 Random Number Generation with NumPy

NumPy random number generators are all stored in the module `numpy.random`. These can be imported with using `import numpy as np` and then calling `np.random.rand`, for example, or by importing `import numpy.random as rnd` and using `rnd.rand.1`.

1.1.1 `rand`, `random_sample`

`rand` and `random_sample` are uniform random number generators which are identical except that `rand` takes a variable number of integer inputs – one for each dimension – while `random_sample` takes a n-element tuple.

`random_sample` is the preferred NumPy function, and `rand` is a convenience function primarily for *MATLAB* users.

```
>>> x = rand(3,4,5)
>>> y = random_sample((3,4,5))
```

1.1.2 `randn`, `standard_normal`

`randn` and `standard_normal` are standard normal random number generators. `randn`, like `rand`, takes a variable number of integer inputs, and `standard_normal` takes an n-element tuple. Both can be called with no arguments to generate a single standard normal (e.g. `randn()`). `standard_normal` is the preferred NumPy function, and `randn` is a convenience function primarily for *MATLAB* users .

```
>>> x = randn(3,4,5)
>>> y = standard_normal((3,4,5))
```

1.1.3 `randint`, `random_integers`

`randint` and `random_integers` are uniform integer random number generators which take 3 inputs, low, high and size. Low is the lower bound of the integers generated, high is the upper and size is a n-element tuple. `randint` and `random_integers` differ in that `randint` generates integers exclusive of the value in high (as do most Python functions), while `random_integers` includes the value in high in its range.

```
>>> x = randint(0,10,(100))
>>> x.max() # Is 9 since range is [0,10)
9
>>> y = random_integers(0,10,(100))
>>> y.max() # Is 10 since range is [0,10]
10
```

1.1.4 shuffle

shuffle randomly reorders the elements of an array in place.

```
>>> x = arange(10)
>>> shuffle(x)
>>> x
array([4, 6, 3, 7, 9, 0, 2, 1, 8, 5])
```

1.1.5 permutation

permutation returns randomly reordered elements of an array as a copy while not directly changing the input.

```
>>> x = arange(10)
>>> permutation(x)
array([2, 5, 3, 0, 6, 1, 9, 8, 4, 7])
>>> x
array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
```

NumPy provides a large selection of random number generators for specific distribution. All take between 0 and 2 required inputs which are parameters of the distribution, plus a tuple containing the size of the output. All random number generators are in the module `numpy.random`.

1.2 Simulation and Random Number Generation

Computer simulated random numbers are usually constructed from very complex but ultimately deterministic functions. Because they are not actually random, simulated random numbers are generally described to as **pseudo-random**.

All pseudo-random numbers in NumPy use one core random number generator based on the ***Mersenne Twister***, a generator which can produce a very long series of pseudo-random data before repeating (up to $2^{19937} - 1$ non-repeating values) .

1.2.1 RandomState

`RandomState` is the class used to control the random number generators. Multiple generators can be initialized by `RandomState`.

```
>>> gen1 = np.random.RandomState()
>>> gen2 = np.random.RandomState()
>>> gen1.uniform() # Generate a uniform
0.6767614077579269
>>> state1 = gen1.get_state()
>>> gen1.uniform()
0.6046087317893271
>>> gen2.uniform() # Different, since gen2 has different seed
0.04519705909244154
>>> gen2.set_state(state1)
>>> gen2.uniform() # Same uniform as gen1 after assigning state
0.6046087317893271
```

1.2.2 State

Pseudo-random number generators track a set of values known as the *state*. The state is usually a vector which has the property that if two instances of the same pseudo-random number generator have the same state, the sequence of pseudo-random numbers generated will be identical. The state of the default random number generator can be read using `numpy.random.get_state` and can be restored using `numpy.random.set_state`.

```
>>> st = get_state()
```

```
>>> randn(4)
array([ 0.37283499, 0.63661908, 1.51588209,
        1.36540624])
>>> set_state(st)
>>> randn(4)
array([ 0.37283499, 0.63661908, 1.51588209,
        1.36540624])
```

The two sequences are identical since the state is the same when `randn` is called. The state is a 5- element tuple where the second element is a 625 by 1 vector of unsigned 32-bit integers. In practice the state should only be stored using `get_state` and restored using `set_state`.

1.2.3 `get_state`

`get_state()` gets the current state of the random number generator, which is a 5-element tuple. It can be called as a function, in which case it gets the state of the default random number generator, or as a method on a particular instance of `RandomState`.

`set_state`

`set_state(state)` sets the state of the random number generator. It can be called as a function, in which case it sets the state of the default random number generator, or as a method on a particular instance of `RandomState`. `set_state` should generally only be called using a state tuple returned by `get_state`.

`seed`

`numpy.random.seed` is a more useful function for initializing the random number generator, and can be used in one of two ways. `seed()` will initialize (or reinitialize) the random number generator using some actual random data provided by the operating system.

`seed(s)` takes a vector of values (can be scalar) to initialize the random number generator at particular state. `seed(s)` is particularly useful for producing simulation studies which are reproducible. In the following example, calls to `seed()` produce different random numbers, since these reinitialize using random data from the computer, while calls to `seed(0)` produce the same (sequence) of random numbers.

```
>>> seed()
>>> randn()
array([ 0.62968838])
>>> seed()
>>> randn()
array([ 2.230155])
>>> seed(0)
>>> randn()
array([ 1.76405235])
>>> seed(0)
>>> randn()
array([ 1.76405235])
```

NumPy always calls `seed()` when the first random number is generated. As a result, calling `standard_normal()` across two “fresh” sessions will not produce the same random number.

1.3 Probability Distributions

1.3.1 normal

`normal()` generates a set of random numbers from a standard Normal (Gaussian). `normal(mu, sigma)` generates draws from a Normal distribution with mean μ and standard deviation σ . `normal(mu, sigma, (10,10))` generates a 10 by 10 array of draws from a Normal with mean μ and standard deviation σ .

`normal(mu, sigma)` is equivalent to `mu + sigma * standard_normal()`.

1.3.2 poisson

`poisson()` generates a set of random numbers from a Poisson distribution with $\lambda = 1$.

`poisson(lambda)` generates a draw from a Poisson distribution with expectation λ . `poisson(lambda, (10,10))` generates a 10 by 10 array of draws from a Poisson distribution with expectation λ .

1.3.3 standard_t

`standard_t(nu)` generates a set of random numbers from a Student's t with shape parameter ν .

`standard_t(nu, (10,10))` generates a 10 by 10 array of draws from a Student's t with shape parameter ν .

1.3.4 uniform

`uniform()` generates a uniform random variable on $(0, 1)$.

`uniform(low, high)` generates a uniform on (l, h) . `uniform(low, high, (10,10))` generates a 10 by 10 array of uniforms on (l, h) .

1.4 Continuous Random Variables

SciPy contains a large number of functions for working with continuous random variables. Each function resides in its own class (e.g. `norm` for Normal or `gamma` for Gamma), and classes expose methods for random number generation, computing the PDF, CDF and inverse CDF, fitting parameters using MLE, and computing various moments. The methods are listed below, where `dist` is a generic placeholder for the distribution name in SciPy.

- **`dist.rvs`**
Pseudo-random number generation. Generically, `rvs` is called using `dist.rvs(*args, loc=0, scale=1, size=size)` where `size` is an n-element tuple containing the size of the array to be generated.
- **`dist.pdf`**
Probability density function evaluation for an array of data (element-by-element). Generically, `pdf` is called using `dist.pdf(x, *args, loc=0, scale=1)` where `x` is an array that contains the values to use when evaluating PDF.
- **`dist.logpdf`**
Log probability density function evaluation for an array of data (element-by-element). Generically, `logpdf` is called using `dist.logpdf(x, *args, loc=0, scale=1)` where `x` is an array that contains the values to use when evaluating log PDF.
- **`dist.cdf`**
Cumulative distribution function evaluation for an array of data (element-by-element). Generically, `cdf` is called using `dist.cdf(x, *args, loc=0, scale=1)` where `x` is an array that contains the values to use when evaluating CDF.
- **`dist.ppf`**
Inverse CDF evaluation (also known as percent point function) for an array of values between 0 and 1. Generically, `ppf` is called using `dist.ppf(p, *args, loc=0, scale=1)` where `p` is an array with all elements between 0 and 1 that contains the values to use when evaluating inverse CDF.
- **`dist.fit`**
Estimate shape, location, and scale parameters from data by maximum likelihood using an array of data. Generically, `fit` is called using `dist.fit(data, *args, floc=0, fscale=1)` where `data` is a data array used to estimate the parameters. `floc` forces the location to a particular value (e.g. `floc=0`). `fscale` similarly forces the scale to a particular value (e.g. `fscale=1`). It is necessary to use `floc` and/or `fscale` when computing MLEs if the distribution does not have a location and/or scale. For example, the gamma distribution

is defined using 2 parameters, often referred to as shape and scale. In order to use ML to estimate parameters from a gamma, `floc=0` must be used.

- **dist.median**
Returns the median of the distribution. Generically, median is called using `dist.median(*args, loc=0, scale=1)`.
- **dist.mean**
Returns the mean of the distribution. Generically, mean is called using `dist.mean(*args, loc=0, scale=1)`.
- **dist.moment**
nth non-central moment evaluation of the distribution. Generically, moment is called using `dist.moment(r, *args, loc=0, scale=1)` where r is the order of the moment to compute.
- **dist.varr**
Returns the variance of the distribution. Generically, var is called using `dist.var(*args, loc=0, scale=1)`.
- **dist.std**
Returns the standard deviation of the distribution. Generically, std is called using `dist.std(*args, loc=0, scale=1)`.

1.4.1 Example

The gamma distribution is used as an example. The gamma distribution takes 1 shape parameter a (a is the only element of `*args`), which is set to 2 in all examples.

```
>>> import scipy.stats as stats
>>> gamma = stats.gamma
>>> gamma.mean(2), gamma.median(2), gamma.std(2), gamma.var(2)
(2.0, 1.6783469900166608, 1.4142135623730951, 2.0)
>>> gamma.moment(2,2) gamma.
moment(1,2)**2 # Variance
>>> gamma.cdf(5, 2), gamma.pdf(5, 2)
(0.95957231800548726, 0.033689734995427337)
>>> gamma.ppf(.95957231800548726, 2)
5.0000000000000018
>>> log(gamma.pdf(5, 2)) gamma.
logpdf(5, 2)
0.0
```

```
>>> gamma.rvs(2, size=(2,2))
array([[ 1.83072394,  2.61422551],
       [ 1.31966169,  2.34600179]])
>>> gamma.fit(gamma.rvs(2, size=(1000)), floc = 0) # a, 0, shape
(2.209958533078413, 0, 0.89187262845460313)
```