

数据计算





数据处理及分析常用库

- Numpy：使用Python进行科学计算的基础软件包
- Pandas：基于Numpy的分析结构化数据的工具集，是数据挖掘和数据分析的核心支持库，也提供数据预处理功能
- Scipy：基于 Numpy 的科学计算库，用于数学、科学、工程学等领域



内容

- NumPy
- Pandas
- Scipy



Numpy

Numerical+Python: Python数组计算、矩阵计算和科学计算的核心库

- C语言实现
- 高性能的数组对象
- 用于数组和矩阵操作的大量函数和方法
- 线性代数、傅里叶变换、随机数生成、图形操作等功能
- 可整合C/C++/Fortran代码



numpy.org



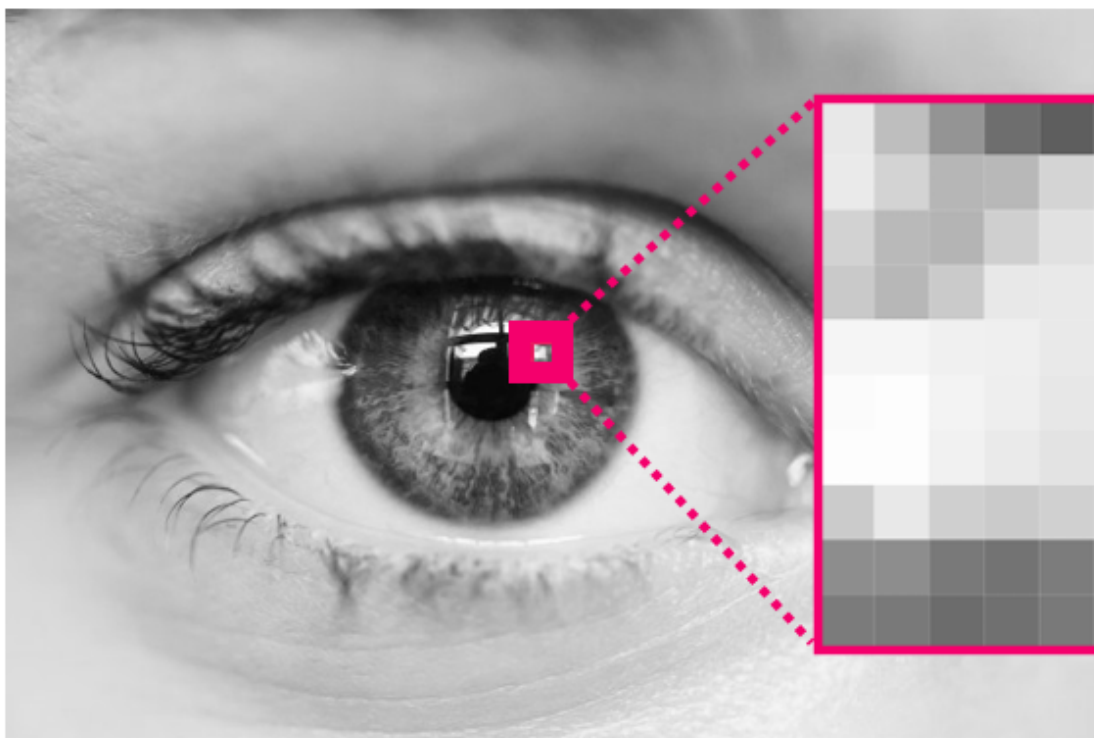
NumPy 中文网

www.numpy.org.cn

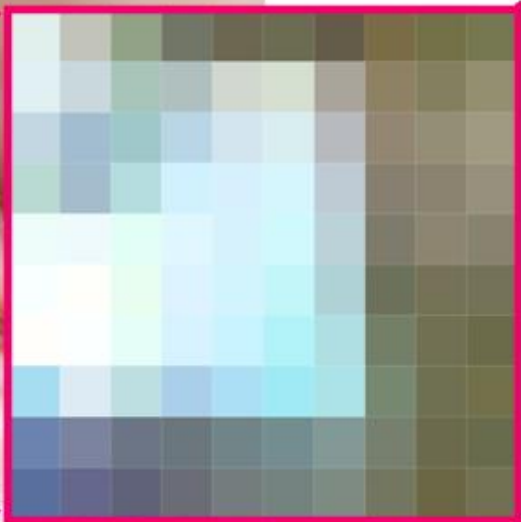
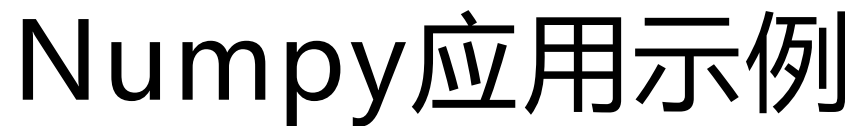
```
import numpy as np
```



Numpy应用示例



230	194	147	108	90	98	84	96	91	101
237	206	188	195	207	213	163	123	116	128
210	183	180	205	224	234	188	122	134	147
198	189	201	227	229	232	200	125	127	135
249	241	237	244	232	226	202	116	125	126
251	254	241	239	230	217	196	102	103	99
243	255	240	231	227	214	203	116	95	91
204	231	208	200	207	201	200	121	95	95
144	140	120	115	125	127	143	118	92	91
121	121	108	109	122	121	134	106	86	97



		233	188	137	96	90	95	63	73	73	82
	237	202	159	120	105	110	88	107	112	121	109
226	191	147	110	101	112	98	123	110	119	142	131
221	191	176	182	203	214	169	144	133	145	155	122
185	160	161	184	205	223	186	137	147	161	140	115
181	174	189	207	206	215	194	136	142	151	133	87
246	237	237	231	208	206	192	122	143	144	111	74
254	254	241	224	199	192	181	99	122	117	107	74
239	248	232	207	187	182	184	110	114	110	113	74
193	215	193	167	158	164	181	114	112	111	105	82
113	119	110	111	113	123	135	120	108	106	113	
93	97	91	103	107	111	122	112	104	114		



NumPy

- 数组的概念
- 数组的创建
- 数组的基本操作
- 矩阵的基本操作
- 常用统计分析函数



数组的相关概念

Description

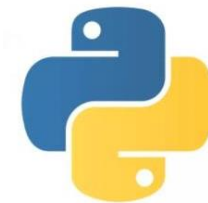
NumPy is the fundamental package for array computing with Python.

Version

1.21.2

Author

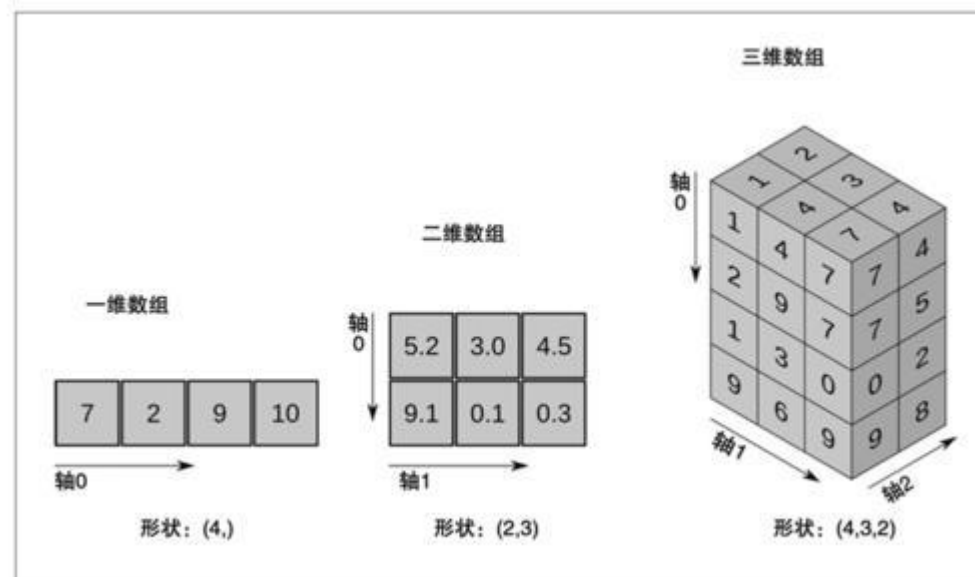
Travis E. Oliphant et al.



数组的概念：数组分类

数组按维度分类

- 一维数组：类似于Python列表
- 二维数组：以一维数组作为数组元素，包括行和列，类似于表格形式，又称矩阵
- 多维数组：维数大于等于三的数组结构。三维数组是最常见的多维数组

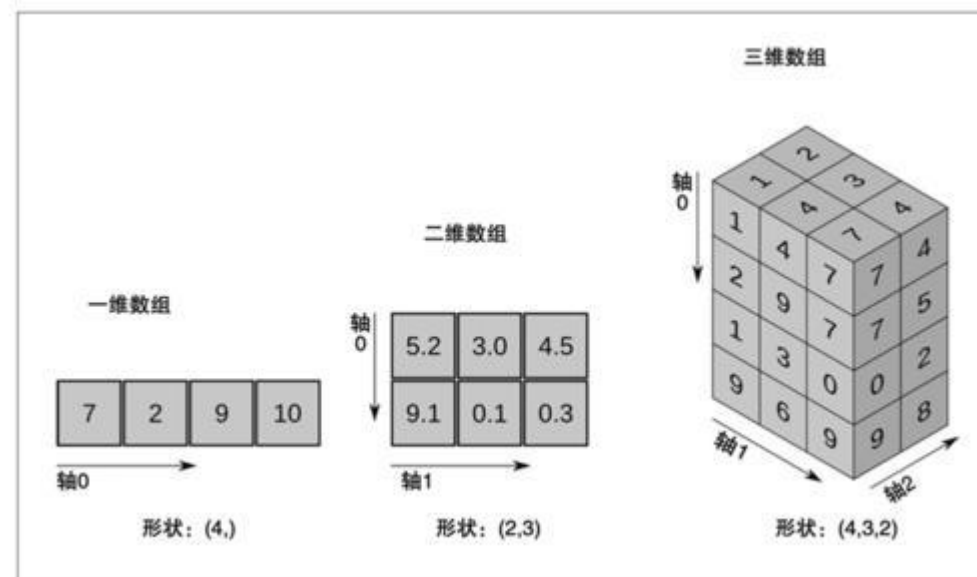




数组的概念：数组的轴

数组的轴 (axis)

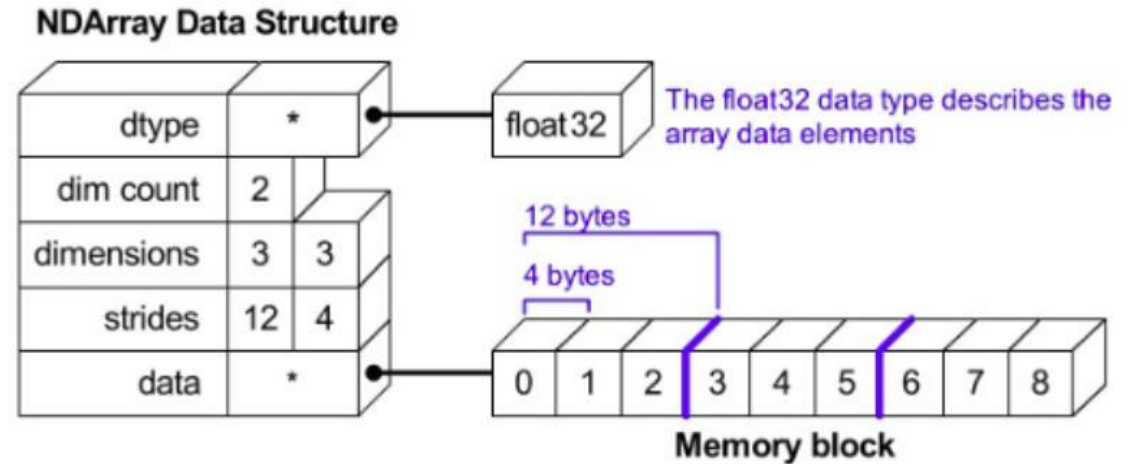
- 数组的维度称为轴 (axis)
- 指定某个axis，就是沿着这个axis做相关操作





NumPy中的数组：ndarray

- 用于存放同类型元素的多维数组
- 每个元素在内存中都有相同大小的存储区域
- 以0下标为开始进行元素的索引



Python View :

0	1	2
3	4	5
6	7	8



数组元素的数据类型

数据类型	描述
bool_	存储一个字节的布尔值 (True, False)
int_	默认整数, 相对于C的long, 通常为int32
intc	相当于C的int, 通常为int32
intp	用于索引的整数, 相当于C的size_t, 通常为int64
int8	字节 (-128~127)
int16	16位整数 (-32768~32767)
int32	32位整数 (-2147483648~2147483647)
int64	64位整数 (-9223372036854775808~9223372036854775807)
uint8	8位无符号整数 (0~255)
uint16	16位无符号整数 (0~65535)



数组元素的数据类型

数据类型	描述
uint32	32位无符号整数 (0~4294967295)
uint64	64位无符号整数 (0~18446744073709551615)
float_	float64的简写
float16	半精度浮点数, 包括: 1 个符号位, 5 个指数位, 10 个尾数位
float32	单精度浮点数, 包括: 1 个符号位, 8 个指数位, 23 个尾数位
float64	双精度浮点数, 包括: 1 个符号位, 11 个指数位, 52 个尾数位
complex_	complex128 类型的简写, 即 128 位复数
complex64	复数, 由两个 32 位浮点数表示 (实数部分和虚数部分)
complex128	复数, 由两个 64 位浮点数表示 (实数部分和虚数部分)
datetime64	日期时间类型
timedelta64	两个时间之间的间隔



数据转换函数

每一种数据类型都有相应的数据转换函数

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

```
>>> np.int8(3.1415)
```

```
3
```

```
>>> np.float64(8)
```

```
8.0
```

```
>>> np.float16(True)
```

```
1.0
```

```
>>> np.bool_(3)
```

```
True
```

```
>>> np.float64(8+j)
```

```
Traceback (most recent call last):  
  File "<input>", line 1, in <module>  
NameError: name 'j' is not defined
```



数组的常用属性

- `a = array([0,1,2,3,4])`
- `b = array([(1.0, 2.0, 3.0), (4.0, 5.0, 6.0)])`
- `c = array([[1,2,3], [4,5,6]], dtype=complex)`

属性名	描述	数组a	数组b	数组c
ndim	维数，等于秩	1	2	2
shape	每个维度上的长度	5,	2,3	2,3
size	数组元素的总个数	5	6	6
dtype	表示数组中元素类型的对象	int32	float64	complex128
itemsize	数组中每个元素的大小(以字节为单位)	4	8	16



数组的属性

```
>>> a = np.array([0, 1, 2, 3, 4])
```

```
▼ a = {ndarray: (5,)} [0 1 2 3 4] ...View as Array
> min = {int32: ()} 0
> max = {int32: ()} 4
> 123 shape = {tuple: 1} (5,)
> dtype = {dtype[int32]: ()} int32
01 size = {int} 5
```




NumPy

- 数组的概念
- 数组的创建
- 数组的基本操作
- 矩阵的基本操作
- 常用统计分析函数



创建数组

创建数组有五种常规机制

1. 从其它Python结构（例如，列表，元组）转换
2. numpy原生数组的创建（例如，`arange`、`ones`、`zeros`等）
3. 从磁盘读取数组，无论是标准格式还是自定义格式
4. 通过使用字符串或缓冲区从原始字节创建数组
5. 使用特殊库函数（例如，`random`）



创建数组：由其它结构创建

```
numpy.array(object, dtype=None, copy=True, order= 'K' ,  
            subok=False, ndmin=0)
```

- object: 任何具有数组接口方法的对象
- dtype: 数据类型
- copy: 若为True, 则复制object对象; 否则, 只有当满足某些要求时, 才会进行复制
- order: 元素在内存中出现的顺序, 值可以为K、A、C、F
- subok: 若为True, 则传递子类; 否则返回的数组将强制为基类数组
- ndmin: 指定生成数组的最小维数



创建数组： 由其它结构创建

```
>>> n1 = np.array([1, 2, 3])
```

```
>  n1 = {ndarray: (3,)} [1 2 3] ...View as Array
```

	0	1	2
0	1	2	3

```
>>> n2 = np.array([1, 2.0], [0, 0], (1+1j, 3.))
```

```
>  n2 = {ndarray: (3, 2)} [[1.+0.j 2.+0.j], [0.+0.j 0.+0.j], [1.+1.j 3.+0.j]] ...View as Array
```

	0	1
0	(1+0j)	(2+0j)
1	0j	0j
2	(1+1j)	(3+0j)



创建数组：由其它结构创建

```
>>> list = [1, 2, 3]
```


```
>>> n3 = np.array(list, dtype=np.float_)
```

```
>  n3 = {ndarray: (3,)} [1. 2. 3.] ...View as Array
```

	0	1	2
0	1.00000	2.00000	3.00000

```
>>> n4 = np.array(list, ndmin=3)
```

```
>  n4 = {ndarray: (1, 1, 3)} [[[1 2 3]]] ...View as Array
```

```
 n4[0]
```

	0	1	2
0	1	2	3

```
>>> n4[0]
```

```
array([[1, 2, 3]])
```

```
>>> n4[0][0]
```

```
array([1, 2, 3])
```

```
>>> n4[0][0][0]
```

```
1
```

```
2
```



创建数组：原生数组创建

创建指定形状、数据类型且未初始化的数组

`numpy.empty(shape, dtype=float, order= 'C')`

- order: 元素在内存中存储的顺序, C: 行优先; F: 列优先

```
>>> n7 = np.empty([2,3],dtype=int)
```

```
> n7 = {ndarray: (2, 3)} [[ 0 1072693248 0], [1073741824 0 1074266112]].
```

	0	1	2
0	0	1072693248	0
1	1073741824	0	1074266112



创建数组：原生数组创建

创建指定维度、数据类型，以0填充的数组

```
numpy.zeros(shape, dtype=float, order= 'C' )
```

```
>>> n8 = np.zeros((3, 3))
```

```
>  n8 = {ndarray: (3, 3)} [[0. 0. 0.], [0. 0. 0.], [0. 0. 0.]]
```

	0	1	2
0	0.00000	0.00000	0.00000
1	0.00000	0.00000	0.00000
2	0.00000	0.00000	0.00000



创建数组：原生数组创建

创建指定维度、数据类型，以1填充的数组

```
numpy.ones(shape, dtype=float, order= 'C' )
```

```
>>> n9 = np.ones((4,2),dtype=int)
```

```
>  n9 = {ndarray: (4, 2)} [[1 1], [1 1], [1 1], [1 1]]
```

0		1
0	1	1
1	1	1
2	1	1
3	1	1



创建数组：原生数组创建

创建指定维度、数据类型，以指定值填充的数组

```
numpy.full(shape, fill_value, dtype=None, order= 'C' *, like=None)
```

```
>>> n10 = np.full((3,3),8)
```

```
>  n10 = {ndarray: (3, 3)} [[8 8 8], [8 8 8], [8 8 8]]
```

0		1	2
0	8	8	8
1	8	8	8
2	8	8	8



创建数组：原生数组创建

通过数据范围创建数组

`numpy.arange([start,] stop[, step], dtype=None)`

```
>>> n11 = np.arange(2, 3, 0.2)
```

```
>  n11 = {ndarray: (5,)} [2. 2.2 2.4 2.6 2.8]
```

	0	1	2	3	4
0	2.00000	2.20000	2.40000	2.60000	2.80000



创建数组：原生数组创建

创建等差数列

`numpy.linspace(start, stop, num=50, endpoint=True, retstep=False, dtype=None)`

- start: 序列的起始值
- stop: 序列的终止值
- num: 要生成的等步长的样本数量，默认值为50
- endpoint: 若为True，数列中包含stop参数的值；否则不包括
- retstep: 若为True，则生成的数组中会显示间距，否则不显示



创建数组：原生数组创建

创建等差数列

```
>>> n12 = np.linspace(1,4,6) >>> n13 = np.linspace(1,4,6,endpoint=False)
```

> `n12` = {ndarray: (6,)} [1. 1.6 2.2 2.8 3.4 4.] > `n13` = {ndarray: (6,)} [1. 1.5 2. 2.5 3. 3.5]

n12							✕	
	0	1	2	3	4	5		
0	1.00000	1.60000	2.20000	2.80000	3.40000	4.00000		

n13							✕	
	0	1	2	3	4	5		
0	1.00000	1.50000	2.00000	2.50000	3.00000	3.50000		







创建数组：原生数组创建

创建等差数列

```
>>> n12 = np.linspace(1,4,6,retstep=True)

>>> n12

(array([1. , 1.6, 2.2, 2.8, 3.4, 4. ]), 0.6)
```

```
▼  n12 = {tuple: 2} (array([1. , 1.6, 2.2, 2.8, 3.4, 4. ]), 0.6)
  >  0 = {ndarray: (6,)} [1. 1.6 2.2 2.8 3.4 4. ]...作为Array查看
    01  1 = {float64: ()} 0.6
    01  __len__ = {int} 2
```



创建数组：原生数组创建

创建等比数列

`numpy.logspace(start, stop, num=50, endpoint=True, base=10.0, dtype=None)`

- start: 序列的起始值为 $\text{base}^{\text{start}}$
- stop: 序列的终止值为 $\text{base}^{\text{stop}}$
- num: 要生成的等比的样本数量，默认值为50
- endpoint: 若为True，数列中包含stop参数的值；否则不包括
- base: 序列范围的基数，默认值为10



创建数组：原生数组创建

创建等比数列

```
>>> n14 = np.logspace(0, 9, 10, base=2, dtype='int')
```

```
>  n14 = {ndarray: (10,)} [ 1  2  4  8 16 32 64 128 256 512]
```

	0	1	2	3	4	5	6	7	8	9
0	1	2	4	8	16	32	64	128	256	512



创建数组：使用random模块生成随机数组

生成0~1之间的随机数组

`numpy.random.rand(d0, d1, d2, ...dn)`

- d0、d1~dn：表示数组的shape，可以为空

```
>>> n15 = np.random.rand(5)
```

```
>>> print(n15)
```

```
[0.02790447 0.43099083 0.76816699 0.41171074 0.18020472]
```

```
>>> n16 = np.random.rand(3, 2)
```

```
>>> print(n16)
```

```
[[0.15065915 0.43655708]
```

```
 [0.99148704 0.89508559]
```

```
 [0.05799208 0.63031163]]
```




创建数组：使用random模块生成随机数组

生成一定范围内的随机整数数组

`numpy.random.randint (low, high=None, size=None)`

- low: 低值（起始值），整数
- high: 高值（终止值），整数
- size: 数组shape，整数或元组，表示一维或多维数组。默认为空，如果为空，则仅返回一个整数



创建数组：使用random模块生成随机数组

生成一定范围内的随机整数数组

```
>>> n21 = np.random.randint(5, size=(5, 2))  
  
>>> n19 = np.random.randint(1, 3, 10)  
>>> print(n19)  
[1 2 2 1 2 2 1 1 1 1]  
  
>>> n20 = np.random.randint(5, 10)  
>>> print(n20)  
9  
  
>>> print(n21)  
[[2 2]  
 [0 3]  
 [1 3]  
 [4 2]  
 [0 3]]
```



创建数组：使用random模块生成随机数组

生成满足标准正态分布的随机数组

`numpy.random.randn(d0, d1, d2, ...dn)` `>>> n18 = np.random.randn(3, 2)`

- `d0`、`d1~dn`：表示数组的shape，可以为空

`>>> n17 = np.random.randn(5)`

`>>> print(n17)`

`[-0.42837384 -1.91183818 1.62630172 -2.19883829 1.38377623]`

`>>> print(n18)`

`[[0.09397946 1.14630998]`

`[-1.32981832 1.02390523]`

`[-0.21326614 0.45270102]]`



创建数组：使用random模块生成随机数组

生成满足正态分布的随机数组

`numpy.random.normal(loc, scale, size)`

- loc: 正态分布的均值
- scale: 正态分布的标准差
- size: 数组的shape



创建数组：使用random模块生成随机数组

生成满足正态分布的随机数组

```
>>> n22 = np.random.normal(0,0.1,10)
>>> print(n22)
[ 0.07774851 -0.25509784  0.17550369 -0.02185514 -0.05450839  0.02559789
-0.17311362  0.1642502   0.1486492  -0.11539874]

>>> n23 = np.random.normal(1,0.3,(3,3))
>>> print(n23)
[[0.83502214 0.66895048 0.93064742]
 [0.74280158 0.53162327 1.02717876]
 [1.00990879 1.14287733 1.06037926]]
```



NumPy

- 数组的概念
- 数组的创建
- 数组的基本操作
- 矩阵的基本操作
- 常用统计分析函数



数组的基本操作：数组运算

NumPy数组不需要编写循环即可对数组数据执行批量运算，称为“矢量化”

- 算术运算和比较运算
- 广播运算
- 函数运算



数组运算： 算术运算和比较运算

```
>>> n1 = np.array([[1,2,3],[4,5,6]])
>>> n2 = np.array([[10,20,30],[40,50,66]])
>>> print(n1+n2)
[[11 22 33]
 [44 55 72]]
```




数组运算： 算术运算和比较运算

```
>>> n1 = np.array([1,2])
```

```
>>> n2 = np.array([3,4])
```

```
>>> print(n1+n2)
```

```
[4 6]
```

```
>>> print(n1-n2)
```

```
[-2 -2]
```

```
>>> print(n1*n2)
```

```
[3 8]
```

```
>>> print(n1/n2)
```

```
[0.33333333 0.5      ]
```

```
>>> print(n1**n2)
```

```
[ 1 16]
```

```
>>> print(n1>n2)
```

```
[False False]
```

```
>>> print(n1<n2)
```

```
[ True  True]
```

```
>>> print(n1==n2)
```

```
[False False]
```

```
>>> print(n1!=n2)
```

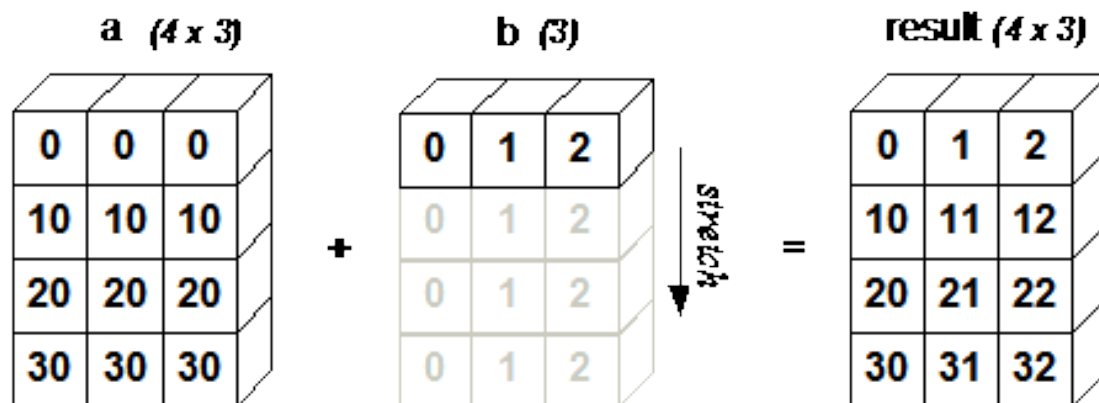
```
[ True  True]
```



数组运算：广播 (Broadcast)

```
a = np.array([[ 0, 0, 0],  
              [10,10,10],  
              [20,20,20],  
              [30,30,30]])  
b = np.array([0,1,2])
```

```
>>> print(a + b)  
[[ 0  1  2]  
 [10 11 12]  
 [20 21 22]  
 [30 31 32]]
```





数组运算：广播（

- 如果两个数组的后缘维度（从末尾开始算起的维度）的轴长度相等或其中一方的长度为1，则认为它们是广播兼容的
- 广播在缺失的维度和（或）轴长度为1的维度上进行

1	2	3		9	9	9	=	.1	.2	.3
4	5	6		9	9	9		.4	.5	.7
7	8	9		9	9	9		.8	.9	1.

1	2	3		-1	0	1	=	-1	0	3
4	5	6		-1	0	1		-4	0	6
7	8	9		-1	0	1		-7	0	9

1	2	3		3	3	3	=	.3	.7	1.
4	5	6		6	6	6		.6	.8	1.
7	8	9		9	9	9		.8	.9	1.

1	2	3		1	1	1	=	1	2	3
1	2	3		2	2	2		2	4	6
1	2	3		3	3	3		3	6	9



数组运算： 函数运算

函数	描述
add()、subtract()、multiply()、divide ()	数组的加减乘除运算
abs()	数组中各元素取绝对值
sqrt()	数组中各元素的平方根
square()	数组中各元素的平方
log()、log10()、log2()	数组中各元素的自然对数和分别以10、2为底的对数
reciprocal()	数组中各元素的倒数
power()	第一个数组中的元素为底数，第二个数组中的元素为幂，进行幂运算
mod()	数组之间相应元素相除后的余数



数组运算：函数运算

函数	描述
<code>around()</code>	数组中各元素指定小数位的四舍五入数
<code>ceil()</code> 、 <code>floor()</code>	数组中各元素向上取整和向下取整
<code>sin()</code> 、 <code>cos()</code> 、 <code>tan()</code>	数组中各元素(角度)的正弦值、余弦值和正切值
<code>modf()</code>	数组中各元素的小数和整数部分分割为两个独立数组
<code>exp()</code>	数组中各元素的指数值
<code>maximum()</code> 、 <code>fmax()</code>	计算数组元素的最大值
<code>minimum()</code> 、 <code>fmin()</code>	计算数组元素的最小值
<code>copysign()</code>	将第二个数组中各个元素的符号赋给第一个数组中对应的元素



数组运算： 函数运算

```
>>> n4 = np.random.randint(-2,8,5)
>>> print(n4)
[ 0  1 -2  4  7]
>>> print(np.abs(n4))
[0 1 2 4 7]
>>> print(np.sqrt(n4))
<input>:1: RuntimeWarning: invalid value encountered in sqrt
[0.          1.          nan  2.          2.64575131]
>>> print(np.exp(n4))
[1.00000000e+00  2.71828183e+00  1.35335283e-01  5.45981500e+01
 1.09663316e+03]
```



数组运算： 函数运算

```
>>> n5 = np.array([[1,2,3],[4,5,6]])
```

```
>>> n6 = np.array([[7,8,9],[1,2,3]])
```

```
>>> print(np.add(n5,n6))
```

```
[[ 8 10 12]
 [ 5  7  9]]
```

```
>>> print(np.subtract(n5,n6))
```

```
[[ -6 -6 -6]
 [  3  3  3]]
```

```
>>> print(np.multiply(n5,n6))
```

```
[[ 7 16 27]
 [ 4 10 18]]
```

```
>>> print(np.divide(n5,n6))
```

```
[[0.14285714 0.25      0.33333333]
 [4.          2.5      2.          ]]
```

```
>>> print(np.power(n5,n6))
```

```
[[  1  256 19683]
 [  4   25  216]]
```



数组的基本操作：索引和切片

NumPy数组元素一般通过索引和切片访问和修改

- NumPy数组使用标准Python语法`array[index]`对数组进行索引
- NumPy数组的切片与Python列表的切片操作类似

`arr[start: stop: step]`



数组的基本操作：索引和切片

```
>>> n8 = np.arange(1,10)
```

```
>>> print(n8)
```

```
[1 2 3 4 5 6 7 8 9]
```

```
>>> print(n8[4])
```

```
5
```

```
>>> print(n8[2:7:2])
```

```
[3 5 7]
```

```
>>> print(n8[3::4])
```

```
[4 8]
```

```
>>> print(n8[3:])
```

```
[4 5 6 7 8 9]
```

```
>>> print(n8[:5])
```

```
[1 2 3 4 5]
```

```
>>> print(n8[::3])
```

```
[1 4 7]
```

```
>>> print(n8[::-1])
```

```
[9 8 7 6 5 4 3 2 1]
```

```
>>> print(n8[:-3:-1])
```

```
[9 8]
```

```
>>> print(n8[-5::-2])
```

```
[5 3 1]
```



数组的基本操作：索引和切片

```
>>> n9 = np.array([[1,2,3],[4,5,6],[7,8,9]])
>>> print(n9)
[[1 2 3]
 [4 5 6]
 [7 8 9]]

>>> print(n9[:2,1:])
[[2 3]
 [5 6]]

>>> print(n9[:, :1])
[[1]
 [4]
 [7]]

>>> print(n9[1,:2])
[4 5]

>>> print(n9[2,1])
8

>>> print(n9[-1])
[7 8 9]

>>> print(n9[:2,2])
[3 6]
```



数组的基本操作：索引和切片

```
>>> n8 = np.arange(1,10)
>>> print(n8)
[1 2 3 4 5 6 7 8 9]
```

```
>>> n9 = n8[2:7:2]
>>> n9
array([3, 5, 7])
```

```
>>> n9[0]=0
>>> n8
array([1, 2, 0, 4, 5, 6, 7, 8, 9])
```

```
>>> n9 = [-1,-2,-3]
>>> n8
array([1, 2, 0, 4, 5, 6, 7, 8, 9])
```



数组的基本操作：布尔索引

```
>>> n9 = np.array([[1,2,3],[4,5,6],[7,8,9]])
```

```
>>> print(n9)
```

```
[[1 2 3]
 [4 5 6]
 [7 8 9]]
```

```
>>> mask = np.array([[True,False,False],[True,False,True],[False,False,True]])
```

```
>>> print(n9[mask])
```

```
[1 4 6 9]
```

```
>>> print(n9[n9>5])
```

```
[6 7 8 9]
```



数组的基本操作：花式索引

```
>>> n1
array([[ 0,  1,  2,  3],
       [ 4,  5,  6,  7],
       [ 8,  9, 10, 11],
       [12, 13, 14, 15],
       [16, 17, 18, 19],
       [20, 21, 22, 23],
       [24, 25, 26, 27],
       [28, 29, 30, 31]])
```

```
>>> n1[[4,2,6]]
array([[16, 17, 18, 19],
       [ 8,  9, 10, 11],
       [24, 25, 26, 27]])
```

```
>>> n1[[-4,-2,-6]]
array([[16, 17, 18, 19],
       [24, 25, 26, 27],
       [ 8,  9, 10, 11]])
```

```
>>> n1[[1,5,7,2],[0,3,2,1]]
array([ 4, 23, 30,  9])
```



数组的基本操作：数组重塑

数组重塑指更改数组的形状 (shape) , 即调整维度。数组重塑不改变数组中元素的数量和类型

- 通用重塑方法
- 数组展平
- 数组转置



数组重塑：通用重塑方法

`array.reshape(shape)`

```
>>> n1 = np.arange(6)
>>> print(n1)
[0 1 2 3 4 5]
```

```
>>> n2 = n1.reshape(2,3)
>>> n3 = n1.reshape(3,2)
>>> n4 = n1.reshape(6,1)
```

```
>>> print(n2)
[[0 1 2]
 [3 4 5]]
```

```
>>> print(n3)
[[0 1]
 [2 3]
 [4 5]]
```

```
>>> print(n4)
[[0]
 [1]
 [2]
 [3]
 [4]
 [5]]
```



数组重塑：通用重塑方法

`array.reshape(shape)`

```
>>> n1 = np.arange(6)
>>> print(n1)
[0 1 2 3 4 5]
```

```
>>> n2 = n1.reshape(2,3)
>>> n3 = n1.reshape(3,2)
>>> n4 = n1.reshape(6,1)
```

```
>>> n1[2] = 10
>>> print(n1)
[ 0  1 10  3  4  5]
```

```
>>> print(n2)
[[ 0  1 10]
 [ 3  4  5]]
```

```
>>> print(n3)
[[ 0  1]
 [10  3]
 [ 4  5]]
```

```
>>> print(n4)
[[ 0]
 [ 1]
 [10]
 [ 3]
 [ 4]
 [ 5]]
```




数组重塑：数组展平

array.flatten()

```
>>> n4= np.array([[n+m*10 for n in range(3)]for m in range(2)])
>>> print(n4)
[[ 0  1  2]
 [10 11 12]]

>>> n5 = n4.flatten()
>>> print(n5)
[ 0  1  2 10 11 12]
```

```
>>> n5[2:4] = -1
>>> print(n5)
[ 0  1 -1 -1 11 12]

>>> print(n4)
[[ 0  1  2]
 [10 11 12]]
```



数组重塑：数组转置

array.T

```
>>> n6 = np.arange(24).reshape(4,6)
```

```
>>> print(n6)
```

```
[[ 0  1  2  3  4  5]
 [ 6  7  8  9 10 11]
 [12 13 14 15 16 17]
 [18 19 20 21 22 23]]
```

```
>>> print(n6.T)
```

```
[[ 0  6 12 18]
 [ 1  7 13 19]
 [ 2  8 14 20]
 [ 3  9 15 21]
 [ 4 10 16 22]
 [ 5 11 17 23]]
```



数组重塑：数组转置

`array.transpose()`

```
>>> n7 = np.array([[ 'A', 100], [ 'B', 200], [ 'C', 300], [ 'D', 400], [ 'E', 500]])
```

```
>>> n8 = n7.transpose()
```



数组重塑：数组转置

PC n7			X	
		0	1	
0	A		100	
1	B		200	
2	C		300	
3	D		400	
4	E		500	

PC n8			X				
		0	1	2	3	4	
0	A	B	C	D	E		
1	100	200	300	400	500		



数组的基本操作：数组元素增加

`numpy.hstack((arr1, arr2,... arrn))`

```
>>> n8 = np.arange(6).reshape((2,3))
>>> n9 = np.arange(10,16).reshape(2,3)
>>> print(n8)
[[0 1 2]
 [3 4 5]]
>>> print(n9)
[[10 11 12]
 [13 14 15]]
```

```
>>> n10 = np.hstack((n8,n9))
>>> print(n10)
[[ 0  1  2 10 11 12]
 [ 3  4  5 13 14 15]]
```

```
>>> n11 = np.hstack((n8,n9,n10))
>>> print(n11)
[[ 0  1  2 10 11 12  0  1  2 10 11 12]
 [ 3  4  5 13 14 15  3  4  5 13 14 15]]
```



数组的基本操作：数组元素增加

`numpy.vstack((arr1, arr2,... arrn))`

```
>>> n8 = np.arange(6).reshape((2,3))
>>> n9 = np.arange(10,16).reshape(2,3)
>>> print(n8)
[[0 1 2]
 [3 4 5]]
>>> print(n9)
[[10 11 12]
 [13 14 15]]
```

```
>>> n12 = np.vstack((n8,n9))
>>> print(n12)
[[ 0  1  2]
 [ 3  4  5]
 [10 11 12]
 [13 14 15]]
```



数组的基本操作：数组元素删除

`numpy.delete(arr, obj, axis=None)`

```
>>> n13 = np.arange(16).reshape(4,4)
>>> print(n13)
[[ 0  1  2  3]
 [ 4  5  6  7]
 [ 8  9 10 11]
 [12 13 14 15]]
```

```
>>> n14 = np.delete(n13,2,axis=0)
>>> print(n14)
[[ 0  1  2  3]
 [ 4  5  6  7]
 [12 13 14 15]]
>>> n15 = np.delete(n13,1,axis=1)
>>> print(n15)
[[ 0  2  3]
 [ 4  6  7]
 [ 8 10 11]
 [12 14 15]]
```



数组的基本操作：数组元素删除

`numpy.delete(arr, obj, axis=None)`

```
>>> n13 = np.arange(16).reshape(4,4)
>>> print(n13)
[[ 0  1  2  3]
 [ 4  5  6  7]
 [ 8  9 10 11]
 [12 13 14 15]]

>>> n16 = np.delete(n13, (1,3), axis=0)
>>> print(n16)
[[ 0  1  2  3]
 [ 8  9 10 11]]

>>> n17 = np.delete(n13, (1,3))
>>> print(n17)
[ 0  2  4  5  6  7  8  9 10 11 12 13 14 15]
```




数组的基本操作：数组元素修改

```
>>> n8 = np.arange(6).reshape((2,3))
```

```
>>> print(n8)
```

```
[[0 1 2]
```

```
 [3 4 5]]
```

```
>>> n8[1]=[10,20,30]
```

```
>>> print(n8)
```

```
[[ 0  1  2]
```

```
 [10 20 30]]
```

```
>>> n8[0][0],n8[1][0] = 5,15
```

```
>>> print(n8)
```

```
[[ 5  1  2]
```

```
 [15 20 30]]
```



数组的基本操作：数组条件查询

`numpy.where(condition, x, y)`

```
>>> n15 = np.arange(9).reshape(3,3)
>>> print(n15)
[[0 1 2]
 [3 4 5]
 [6 7 8]]
```

```
>>> print(np.where(n15>5,1,-1))
[[-1 -1 -1]
 [-1 -1 -1]
 [ 1  1  1]]
```



NumPy

- 数组的概念
- 数组的创建
- 数组的基本操作
- 矩阵的基本操作
- 常用统计分析函数



NumPy中的数组（Array）和矩阵（Matrix）

- 矩阵即二维数组，是数组的分支
- Array和Matrix是NumPy中的两种不同的数据类型
- 在Array上可以执行的运算及函数，一般都可用于Matrix
- 在Array和Matrix上执行相同的数学运算时，可能得到不同的结果



矩阵创建

`numpy.mat(data, dtype=None)`

```
>>> m1 = np.mat([[1, 2], [3, 4]])
```

```
>>> print(m1)
```

```
[[1 2]
 [3 4]]
```

```
>>> m1
```

```
matrix([[1, 2],
        [3, 4]])
```



矩阵创建

`numpy.mat(data, dtype=None)`

```
>>> m1 = np.mat('1 2 3;4 5 6')
```

```
>>> print(m1)
```

```
[[1 2 3]
```

```
 [4 5 6]]
```



矩阵创建

`numpy.mat(data, dtype=None)`

```
>>> m2 = np.mat(np.zeros((3,3)))
```

```
>>> print(m2)
```

```
[[0. 0. 0.]
```

```
 [0. 0. 0.]
```

```
 [0. 0. 0.]]
```

```
>>> m3 = np.mat(np.random.randint(1,8,size=(3,5)))
```

```
>>> print(m3)
```

```
[[1 7 7 4 2]
```

```
 [7 7 3 5 5]
```

```
 [2 7 6 1 3]]
```



矩阵运算：+，-，/

```
>>> m4 = np.mat(np.arange(6).reshape((2,3)))
>>> m5 = np.mat(np.arange(10,16).reshape((2,3)))
>>> print(m4)
[[0 1 2]
 [3 4 5]]
>>> print(m5)
[[10 11 12]
 [13 14 15]]

>>> print(m4+m5)
[[10 12 14]
 [16 18 20]]

>>> print(m4-m5)
[[-10 -10 -10]
 [-10 -10 -10]]

>>> print(m4/m5)
[[0.          0.09090909 0.16666667]
 [0.23076923 0.28571429 0.33333333]]
```




矩阵运算：*

```
>>> print(m4)
```

```
[[0 1 2]
```

```
 [3 4 5]]
```

```
>>> print(m5)
```

```
[[10 11 12]
```

```
 [13 14 15]]
```

```
>>> print(m4*m5)
```

```
Traceback (most recent call last):
```

```
File "<input>", line 1, in <module>
```

```
File "C:\Users\wangj\PycharmProjects\dataAnalysisProject\venv\lib\site-packages\numpy\matrixlib\defmatrix.py", line 218, in __mul__
```

```
    return N.dot(self, asmatrix(other))
```

```
File "<__array_function__ internals>", line 5, in dot
```



矩阵运算：矩阵相乘

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1s} \\ \cdots & \cdots & \cdots & \cdots \\ a_{i1} & a_{i2} & \cdots & a_{is} \\ \cdots & \cdots & \cdots & \cdots \\ a_{m1} & a_{m2} & \cdots & a_{ms} \end{bmatrix} \cdot \begin{bmatrix} b_{11} & \cdots & b_{1j} & \cdots & b_{1n} \\ b_{21} & \cdots & b_{2j} & \cdots & b_{2n} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ b_{s1} & \cdots & b_{sj} & \cdots & b_{sn} \end{bmatrix}$$

$$= \begin{bmatrix} c_{11} & \cdots & a_{i1}b_{1j} + a_{i2}b_{2j} + \cdots + a_{is}b_{sj} & \cdots & c_{1n} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ c_{i1} & \cdots & c_{ij} & \cdots & c_{in} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ c_{m1} & \cdots & c_{mj} & \cdots & c_{mn} \end{bmatrix}$$

矩阵相乘使用点乘（点积），即对应位置——相乘后求和



矩阵运算：矩阵相乘

```
>>> print(m4)
[[0 1 2]
 [3 4 5]]
>>> print(m5)
[[10 11 12]
 [13 14 15]]
```

```
>>> m5 = m5.T
>>> print(m5)
[[10 13]
 [11 14]
 [12 15]]
```

```
>>> print(m4*m5)
[[ 35  44]
 [134 170]]
```



NumPy数组的点乘

`numpy.dot(arr1, arr2)`

```
>>> n1 = np.array([[1,2],[3,4]])
```

```
>>> n2 = np.array([[5,6],[7,8]])
```

```
>>> print(n1)
```

```
[[1 2]
```

```
 [3 4]]
```

```
>>> print(n2)
```

```
[[5 6]
```

```
 [7 8]]
```

```
>>> print(n1*n2)
```

```
[[ 5 12]
```

```
 [21 32]]
```

```
>>> print(np.dot(n1,n2))
```

```
[[19 22]
```

```
 [43 50]]
```



矩阵求逆

matrix.I

```
>>> m6 = np.mat('1 3 3;4 5 6;7 12 9')
```

```
>>> print(m6)
```

```
[[ 1  3  3]
 [ 4  5  6]
 [ 7 12  9]]
```

```
>>> print(m6.I)
```

```
[[ -0.9          0.3          0.1          ]
 [  0.2          -0.4          0.2          ]
 [  0.43333333  0.3         -0.23333333]]
```



NumPy

- 数组的概念
- 数组的创建
- 数组的基本操作
- 矩阵的基本操作
- 常用统计分析函数



统计分析函数

函数	描述
sum()	对数组中的元素或某行某列的元素求和
cumsum()	数组元素累计求和
cumprod()	数组元素累计求积
mean()	数组元素求平均值
min()、max()	数组的最小值和最大值
average()	计算加权平均数
median()	数组元素的中位数（中值）
var()	方差
std()	标准差



统计分析函数

函数	描述
<code>argmin()</code> 、 <code>argmax()</code>	数组最小值和最大值的下标（默认情况下为一维下标值）
<code>ptp()</code>	计算最大值和最小值的差
<code>np.unravel_index()</code>	根据数组形状将一维下标转换为多维下标
<code>np.percentile()</code>	求百分位数的值



统计分析函数

```
>>> n7 = np.arange(12).reshape(3,4)
>>> n7
array([[ 0,  1,  2,  3],
       [ 4,  5,  6,  7],
       [ 8,  9, 10, 11]])
```

```
>>> n7.sum()
66
>>> n7.sum(axis=0)
array([12, 15, 18, 21])
>>> n7.sum(axis=1)
array([ 6, 22, 38])
```

```
>>> n7.cumsum(axis=0)
array([[ 0,  1,  2,  3],
       [ 4,  6,  8, 10],
       [12, 15, 18, 21]], dtype=int32)
```

```
>>> n7.cumsum(axis=1)
array([[ 0,  1,  3,  6],
       [ 4,  9, 15, 22],
       [ 8, 17, 27, 38]], dtype=int32)
```



统计分析函数

```
>>> n7 = np.arange(12).reshape(3,4)
```

```
>>> n7
```

```
array([[ 0,  1,  2,  3],
       [ 4,  5,  6,  7],
       [ 8,  9, 10, 11]])
```

```
>>> n7.argmin()
```

```
0
```

```
>>> n7.argmax()
```

```
11
```

```
>>> np.unravel_index(n7.argmax(), n7.shape)
```

```
(2, 3)
```

```
>>> np.unravel_index((3,9), n7.shape)
```

```
(array([0, 2], dtype=int64), array([3, 1], dtype=int64))
```

```
>>> np.unravel_index([3,9], n7.shape, order='F')
```

```
(array([0, 0], dtype=int64), array([1, 3], dtype=int64))
```



统计分析函数

```
>>> n7 = np.arange(12).reshape(3,4)
>>> n7
array([[ 0,  1,  2,  3],
       [ 4,  5,  6,  7],
       [ 8,  9, 10, 11]])
```

```
>>> np.percentile(n7,50)
5.5
>>> np.percentile(n7,50,axis=0)
array([4., 5., 6., 7.])
```

```
>>> np.percentile(n7,50,axis=1)
array([1.5, 5.5, 9.5])
```

```
>>> np.percentile(n7,30,axis=1)
array([0.9, 4.9, 8.9])
```

```
>>> np.percentile(n7,30,axis=0)
array([2.4, 3.4, 4.4, 5.4])
```

```
>>> np.percentile(n7,90,axis=0)
array([ 7.2,  8.2,  9.2, 10.2])
```

```
>>> np.percentile(n7,90,axis=1)
array([ 2.7,  6.7, 10.7])
```



统计分析函数

```
>>> n7 = np.array([[1,2,3],[4,5,6],[9,10,11]])
```

```
>>> print(np.sum(n7))
```

```
51
```

```
>>> print(np.amax(n7),np.amin(n7))
```

```
11 1
```

```
>>> print(np.median(n7),np.mean(n7))
```

```
5.0 5.666666666666667
```

```
>>> print(np.std(n7),np.var(n7))
```

```
3.39934634239519 11.555555555555557
```

```
>>> print(np.amax(n7,axis=0))
```

```
[ 9 10 11]
```

```
>>> print(np.amin(n7,axis=1))
```

```
[1 4 9]
```

```
>>> print(np.median(n7,axis=0))
```

```
[4. 5. 6.]
```

```
>>> print(np.mean(n7,axis=1))
```

```
[ 2.  5. 10.]
```



统计分析函数

```
>>> score = np.array([92, 78, 95, 86, 100])
```

```
>>> credit = np.array([3, 2, 2, 3, 2])
```

```
>>> print(np.average(score, weights=credit))
```

```
90.0
```



排序函数

```
>>> n3 = np.array([[4,5,3],[6,2,8],[7,1,9]])
```

```
>>> print(n3)
```

```
[[4 5 3]
 [6 2 8]
 [7 1 9]]
```

```
>>> print(np.sort(n3))
```

```
[[3 4 5]
 [2 6 8]
 [1 7 9]]
```

```
>>> print(np.sort(n3,axis=0))
```

```
[[4 1 3]
 [6 2 8]
 [7 5 9]]
```

```
>>> print(np.sort(n3,axis=1))
```

```
[[3 4 5]
 [2 6 8]
 [1 7 9]]
```



内容

- NumPy
- Pandas
- Scipy



Pandas

Pandas: panel data + Python data analysis: 是一款快速、强大、灵活、易用的开源数据分析和操作工具

- 基于Numpy, 是Python的核心数据分析支持库
- 提供了快速、灵活、明确的数据结构
- 可以从各种文件格式如 CSV、JSON、SQL、Microsoft Excel 导入数据
- 可以对各种数据进行运算操作, 并具有数据清洗和数据加工能力

```
import pandas as pd
```



<https://pandas.pydata.org>



Pandas能够处理的数据类型

- 与SQL或Excel表类似的数据
- 有序和无序（非固定频率）的时间序列数据
- 带行列标签的矩阵数据
- 任意其他形式的观测、统计数据集



Pandas主要数据结构

- Series类型
- DataFrame类型



Series类型：带有标签的一维同构数组

由一组数据及与之相关的数据索引（标签）组成

```
>>> print(s)
```

```
> s = {Series: (3,)} (0, 89) (1, 92) (2, 75) ...View as Series
```

```
0      89
```

```
1      9
```

```
0
```

```
89
```

```
2      7
```

```
1
```

```
92
```

```
2
```

```
75
```

```
dtype: int64
```



Series对象：创建

```
s=pandas.Series(data=None, index=None, dtype=None, name=None, copy=False)
```

- data: 数据，支持列表、字典、numpy数组、标量值
- index: 行标签（索引）
- dtype: 数据类型
- name: 所创建Series的名字
- copy: 是否创建输入数据的副本



Series对象创建：由列表创建

```
>>> s = pd.Series([89, 92, 75])
```

```
>>> print(s)
```

隐式索引

0	89
1	92
2	75

```
dtype: int64
```



Series对象创建：由列表创建

```
>>> s1 = pd.Series([89, 92, 75], index=['A', 'B', 'C'])
```

```
>>> print(s1)
```

A	89
B	92
C	75

```
dtype: int64
```

```
> s1 = {Series: (3,)} ('A', 89) ('B', 92) ('C', 75).
```

显式索引



Series对象创建：由ndarray创建

```
>>> import numpy as np
>>> s2 = pd.Series(np.arange(3), index=['a', 'b', 'c'])

>>> print(s2)
a      0
b      1
c      2
dtype: int32
```



Series对象创建：由字典创建

```
>>> dic = {"A":1, "B":2, "C":3, "D":6}
```

```
>>> s3 = pd.Series(dic)
```

```
>>> print(s3)
```

```
A    1
```

```
B    2
```

```
C    3
```

```
D    6
```

```
dtype: int64
```




Series对象创建：由字典创建

```
>>> dic = {"A":1,"B":2,"C":3,"D":6}
```

```
>>> s3 = pd.Series(dic)
```

```
>>> print(s3)
```

A	1
B	2
C	3
D	6

dtype: int64

```
>>> index4 = ['A','B','C','D','E']
```

```
>>> s4 = pd.Series(dic,index4)
```

```
>>> print(s4)
```

A	1.0
B	2.0
C	3.0
D	6.0
E	NaN

dtype: float64

```
>>> pd.isnull(s4)
```

A	False
B	False
C	False
D	False
E	True

dtype: bool



Series对象创建：由Series创建

```
>>> s3
A      1
B      2
C      3
D      6
dtype: int64
```

```
>>> s4 = s3.reindex(['D', 'C', 'E', 'A', 'B'])
>>> s4
D      6.0
C      3.0
E      NaN
A      1.0
B      2.0
dtype: float64
```



Series对象创建：由Series创建

```
>>> s3
```

```
A      1
```

```
B      2
```

```
C      3
```

```
D      6
```

```
dtype: int64
```

```
>>> s4 = s3.reindex(  
...     ['D','C','E','A','B','F'],fill_value=88)
```

```
>>> s4
```

```
D      6
```

```
C      3
```

```
E     88
```

```
A      1
```

```
B      2
```

```
F     88
```

```
dtype: int64
```



Series对象的索引可重复

```
>>> s4 = pd.Series(np.arange(4), index=[1, 2, 4, 1])
```

```
>>> print(s4)
```

```
1      0
```

```
2      1
```

```
4      2
```

```
1      3
```

```
dtype: int32
```



Series对象的属性和方法

pandas.Series — pandas 1.3.3 documentation (pydata.org)

<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.Series.html>



Series对象的常见运算和操作

- Python列表操作
- Numpy一维数组运算及操作
- Python字典操作
- 特有操作



Series对象查看

```
>>> s4.describe()
```

```
count      4.000000
```

```
mean       3.000000
```

```
std        2.160247
```

```
min        1.000000
```

```
25%        1.750000
```

```
50%        2.500000
```

```
75%        3.750000
```

```
max        6.000000
```

```
dtype: float64
```

```
>>> s4.values
```

```
array([ 1.,  2.,  3.,  6., nan])
```

```
>>> s4.index
```

```
Index(['A', 'B', 'C', 'D', 'E'], dtype='object')
```



Series对象的索引和切片：位置索引

```
>>> print(s3)
```

```
A    1
```

```
B    2
```

```
C    3
```

```
D    6
```

```
dtype: int64
```

```
>>> print(s3[0])
```

```
1
```

```
>>> print(s3[-1])
```

```
6
```

```
>>> print(s3[1:3])
```

```
B    2
```

```
C    3
```

```
dtype: int64
```




Series对象的索引和切片：标签索引

```
>>> print(s3)
```

```
A    1
```

```
B    2
```

```
C    3
```

```
D    6
```

```
dtype: int64
```

```
>>> print(s3['A'])
```

```
1
```

```
>>> print(s3['A':'C'])
```

```
A    1
```

```
B    2
```

```
C    3
```

```
dtype: int64
```



Series对象的索引和切片： 标签索引

```
>>> s4 = pd.Series(np.arange(4), index=[1, 2, 4, 1])
>>> print(s4)
1      0
2      1
4      2
1      3
dtype: int32
```

```
>>> print(s4[1])
1      0
1      3
dtype: int32
```



Series对象的索引和切片：标签索引

```
>>> s2 = pd.Series(np.arange(3), index=[5, 6, 7])
```

```
>>> print(s2)
```

```
5    0
```

```
6    1
```

```
7    2
```

```
dtype: int32
```

```
>>> print(s2[6])
```

```
1
```

```
>>> print(s2[5:7])
```

```
Series([], dtype: int32)
```

```
>>> print(s2[0:3])
```

```
5    0
```

```
6    1
```

```
7    2
```

```
dtype: int32
```



Series对象的索引和切片：点索引

```
>>> print(s1)
```

```
A      89
```

```
B      92
```

```
C      75
```

```
dtype: int64
```

```
>>> print(s1.A)
```

```
89
```



Series对象的算术运算

```
>>> s = pd.Series([1,2,3,4,5],index=['a','b','c','d','e'])
```

```
>>> s[1:]+s[:-1]
```

```
a    NaN
b    4.0
c    6.0
d    8.0
e    NaN
dtype: float64
```

s[1:]

索引	数据
b	2
c	3
d	4
e	100

s[:-1]

索引	数据
a	1
b	2
c	3
d	4

+

=

s[1:]+s[:-1]

索引	数据
a	NaN
b	4.0
c	6.0
d	8.0
e	NaN