

Lec 06 - Advanced indexing & Broadcasting

Statistical Computing and Computation

Sta 663 | Spring 2022

Dr. Colin Rundel

NumPy - Advanced Indexing

From last time: subsetting with tuples

Unlike lists, an ndarray can be subset by a tuple containing integers,

```
x = np.arange(6)
x

## array([0, 1, 2, 3, 4, 5])

x[(0,1,3),]

## array([0, 1, 3])

x[(0,1,3)]

## Error in py_call_impl(callable, dots$args, dots
##
## Detailed traceback:
##   File "<string>", line 1, in <module>
```

```
x[0,1,3]

## Error in py_call_impl(callable, dots$args, dots
More next time on why x[(0,1,3)] does not work.
## Detailed traceback:
```

```
x = np.arange(16).reshape((4,4))
x

## array([[ 0,  1,  2,  3],
##        [ 4,  5,  6,  7],
##        [ 8,  9, 10, 11],
##        [12, 13, 14, 15]])

x[(0,1,3), :]

## array([[ 0,  1,  2,  3],
##        [ 4,  5,  6,  7],
##        [12, 13, 14, 15]])

x[:, (0,1,3)]

## array([[ 0,  1,  3],
##        [ 4,  5,  7],
##        [ 8,  9, 11],
##        [12, 13, 15]])
```

Integer array subsetting (lists)

Lists of integers can be used to subset in the same way:

```
x = np.arange(6)
x

## array([0, 1, 2, 3, 4, 5])

x[[0,1,3],]

## array([0, 1, 3])

x[[0,1,3]]

## array([0, 1, 3])

x[[1.,3.]]

## Error in py_call_impl(callable, dots$args, dots$keywords): In
##
## Detailed traceback:
##   File "<string>", line 1, in <module>
```

Note that the `,` is now optional

```
x = np.arange(16).reshape((4,4))
x

## array([[ 0,  1,  2,  3],
##        [ 4,  5,  6,  7],
##        [ 8,  9, 10, 11],
##        [12, 13, 14, 15]])

x[[1,3]]

## array([[ 4,  5,  6,  7],
##        [12, 13, 14, 15]])

x[[1,3], ]

## array([[ 4,  5,  6,  7],
##        [12, 13, 14, 15]])

x[:, [1,3]]

## array([[ 1,  3],
##        [ 5,  7],
##        [ 9, 11],
##        [13, 15]])

x[[1,3], [1,3]]
```

Integer array subsetting (ndarrays)

Similarly we can also use integer ndarrays:

```
x = np.arange(6)
y = np.array([0,1,3])
z = np.array([1., 3.])

x[y,]

## array([0, 1, 3])

x[y]

## array([0, 1, 3])

x[z]

## Error in py_call_impl(callable, dots$args, dots$keywords): In
##
## Detailed traceback:
## File "<string>", line 1, in <module>
```

Again the , is now optional

```
x = np.arange(16).reshape((4,4))
x

## array([[ 0,  1,  2,  3],
##        [ 4,  5,  6,  7],
##        [ 8,  9, 10, 11],
##        [12, 13, 14, 15]])

y = np.array([1,3])

x[y]

## array([[ 4,  5,  6,  7],
##        [12, 13, 14, 15]])

x[y, ]

## array([[ 4,  5,  6,  7],
##        [12, 13, 14, 15]])

x[:, y]

## array([[ 1,  3],
##        [ 5,  7],
##        [ 9, 11],
##        [13, 15]])
```

Exercise 1

Given the following matrix,

```
x = np.arange(16).reshape((4,4))  
x
```

```
## array([[ 0,  1,  2,  3],  
##        [ 4,  5,  6,  7],  
##        [ 8,  9, 10, 11],  
##        [12, 13, 14, 15]])
```

write an expression to obtain the center 2x2 values (i.e. 5, 6, 9, 10 as a matrix).

Boolean indexing

Lists or ndarrays of boolean values can also be used to subset, positions with `True` are kept and `False` are discarded.

```
x = np.arange(6)
x
```

```
## array([0, 1, 2, 3, 4, 5])
```

```
x[[True, False, True, False, True, False]]
```

```
## array([0, 2, 4])
```

```
x[np.array([True, True, False, False, True, False])]
```

```
## array([0, 1, 4])
```

the utility comes from vectorized comparison operations,

```
x > 3
```

```
## array([False, False, False, False,  True,  True])
```

```
x[x>3]
```

```
## array([4, 5])
```

```
x % 2 == 1
```

```
y = np.arange(9).reshape((3,3))
y % 2 == 0
```

```
## array([[ True, False,  True],
##        [False,  True, False],
##        [ True, False,  True]])
```

```
y[y % 2 == 0]
```

NumPy and Boolean operators

If we want to use a boolean operator on an array we need to use `&`, `|`, and `~` instead of `and`, `or`, and `not` respectively.

```
x = np.arange(6)
x
```

```
## array([0, 1, 2, 3, 4, 5])
```

```
y = x % 2 == 0
y
```

```
## array([ True, False,  True, False,  True, False])
```

```
~y
```

```
## array([False,  True, False,  True, False,  True])
```

```
y & (x > 3)
```

```
## array([False, False, False, False,  True, False])
```

```
y | (x > 3)
```


meshgrid

One other useful function in NumPy is `meshgrid()` which generates all possible combinations between the input vectors,

```
pts = np.arange(3)
x, y = np.meshgrid(pts, pts)
x
```

```
## array([[0, 1, 2],
##        [0, 1, 2],
##        [0, 1, 2]])
```

```
y
```

```
## array([[0, 0, 0],
##        [1, 1, 1],
##        [2, 2, 2]])
```

```
np.sqrt(x**2 + y**2)
```

```
## array([[0.          , 1.          , 2.          ],
##        [1.41421356, 2.23606798],
##        [2.23606798, 2.82842712]])
```

Exercise 2

We will now use this to attempt a simple brute force approach to numerical optimization, define a grid of points using `meshgrid()` to approximate the minima the following function:

$$f(x, y) = (1 - x)^2 + 100(y - x^2)^2$$

Considering values of $x, y \in (-1, 3)$, which values of x, y minimize this function?

NumPy - Broadcasting

Broadcasting

This is an approach for deciding how to generalize arithmetic operations between arrays with differing shapes.

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

```
x * 2
```

```
## array([2, 4, 6])
```

```
x * np.array([2])
```

```
## array([2, 4, 6])
```

```
x * np.array([2,2,2])
```

```
## array([2, 4, 6])
```

In the first example 2 is equivalent to the array `np.array([2])` which is being broadcast across the longer array `x`.

Efficiency

Using broadcasts can be much more efficient as it does not copy the underlying data,

```
x = np.arange(1e5)  
y = np.array([2]).repeat(1e5)
```

```
%timeit x * 2
```

31.3 μs \pm 1.3 μs per loop (mean \pm std. dev. of 7 runs, 10,000 loops each)

```
%timeit x * y
```

70.5 μs \pm 2.93 μs per loop (mean \pm std. dev. of 7 runs, 10,000 loops each)

General Broadcasting

When operating on two arrays, NumPy compares their shapes element-wise. It starts with the trailing (i.e. rightmost) dimensions and works its way left. Two dimensions are compatible when

1. they are equal, or
2. one of them is 1

If these conditions are not met, a `ValueError: operands could not be broadcast together` exception is thrown, indicating that the arrays have incompatible shapes. The size of the resulting array is the size that is not 1 along each axis of the inputs.

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

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

```
x + np.array([1,2,3])
```

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

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

```
x + np.array([1,2,3])
```

A quick fix

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

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

```
x + np.array([1,2,3])
```

```
## Error in py_call_impl(callable, dots$args, dots$keywords): ValueError: operands could not be broadcast together with shapes (3,4),(1,3)  
##  
## Detailed traceback:  
##   File "<string>", line 1, in <module>
```

```
x + np.array([1,2,3]).reshape(3,1)
```

```
## array([[ 1,  2,  3,  4],  
##        [ 6,  7,  8,  9],  
##        [11, 12, 13, 14]])
```


Mechanics

```
x = np.arange(12).reshape((4,3))
y = 1
x+y
```

```
## array([[ 1,  2,  3],
##        [ 4,  5,  6],
##        [ 7,  8,  9],
##        [10, 11, 12]])
```

```
x    (2d array): 4 x 3
y    (1d array):    1
-----
x+y  (2d array): 4 x 3
```

```
x = np.arange(12).reshape((4,3))
y = np.array([1,2,3])
x+y
```

```
## array([[ 1,  3,  5],
##        [ 4,  6,  8],
##        [ 7,  9, 11],
##        [10, 12, 14]])
```

```
x = np.arange(12).reshape((3,4))
y = np.array([1,2,3])
x+y
```

```
## Error in py_call_impl(callable, dots$args, dots$keys)
##
## Detailed traceback:
##   File "<string>", line 1, in <module>
```

```
x    (2d array): 3 x 4
y    (1d array):    3
-----
x+y  (2d array): Error
```

```
x = np.arange(12).reshape((3,4))
y = np.array([1,2,3]).reshape((3,1))
x+y
```

```
## array([[ 1,  2,  3,  4],
##        [ 6,  7,  8,  9],
##        [11, 12, 13, 14]])
```

Another example

```
a = np.array([0,10,20,30]).reshape((4,1))  
b = np.array([1,2,3])
```

a

```
## array([[ 0],  
##       [10],  
##       [20],  
##       [30]])
```

b

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

a+b

```
## array([[ 1,  2,  3],  
##       [11, 12, 13],  
##       [21, 22, 23],  
##       [31, 32, 33]])
```

a (2d array): 4 x 1

b (1d array): 3

x+y (2d array): 4 x 3



From NumPy user guide - [Broadcasting](#)

Example - Standardizing

Below we generate a data set with 3 columns of random normal values. Each column has a different mean and standard deviation which we can check with `mean()` and `std()`.

```
rng = np.random.default_rng(1234)

d = rng.normal(loc=[-1,0,1], scale=[1,2,3], size=(1000,3))
d.mean(axis=0)
```

```
## array([-1.0294382 , -0.01396257,  1.01241784])
```

```
d.std(axis=0)
```

```
## array([0.99674719,  2.03222595,  3.10625219])
```

Use broadcasting to standardize all three columns to have mean 0 and standard deviation 1.

Check the new data set using `mean()` and `std()`.

Exercise 3

For each of the following combinations determine what the resulting dimension will be:

- $A (128 \times 128 \times 3) + B (3)$
- $A (8 \times 1 \times 6 \times 1) + B (7 \times 1 \times 5)$
- $A (2 \times 1) + B (8 \times 4 \times 3)$
- $A (3 \times 1) + B (15 \times 3 \times 5)$
- $A (3) + B (4)$

Broadcasting and assignment

In addition to arithmetic operators, broadcasting can be used with assignment via array indexing,

```
x = np.arange(12).reshape((3,4))
y = -np.arange(4)
z = -np.arange(3)
```

```
x[:] = y
x
```

```
## array([[ 0, -1, -2, -3],
##        [ 0, -1, -2, -3],
##        [ 0, -1, -2, -3]])
```

```
x[...] = y
x
```

```
## array([[ 0, -1, -2, -3],
##        [ 0, -1, -2, -3],
##        [ 0, -1, -2, -3]])
```

```
x[:] = z
```

```
## Error in py_call_impl(callable, dots$args, dots$keywords)
##
## Detailed traceback:
##   File "<string>", line 1, in <module>
```

```
x[:] = z.reshape((3,1))
x
```

```
## array([[ 0,  0,  0,  0],
##        [-1, -1, -1, -1],
##        [-2, -2, -2, -2]])
```

NumPy - Basic file IO

Reading and writing arrays

We will not spend much time on this as most data you will encounter is more likely to be a tabular format (e.g. data frame) and tools like Pandas are more appropriate.

For basic saving and loading of NumPy arrays there are the `save()` and `load()` functions which use a built in binary format.

```
x = np.arange(1e5)
np.save("data/x.npy", x)
new_x = np.load("data/x.npy")
np.all(x == new_x)
```

```
## True
```

Additional functions for saving (`savez()`, `savez_compressed()`, `savetxt()`) exist for saving multiple arrays or saving a text representation of an array.

Reading delimited data

While not particularly recommended, if you need to read delimited (csv, tsv, etc.) data into a NumPy array you can use `genfromtxt()`,

```
options(width=300)
```

```
with open("data/mtcars.csv") as file:
    mtcars = np.genfromtxt(file, delimiter=";", skip_header=True)
```

```
mtcars
```

```
## array([[ 6.  , 160.  , 110.  ,  3.9  ,  2.62 , 16.46 ,  0.  ,  1.  ,  4.  ,  4.  ],
##        [ 6.  , 160.  , 110.  ,  3.9  ,  2.875, 17.02 ,  0.  ,  1.  ,  4.  ,  4.  ],
##        [ 4.  , 108.  ,  93.  ,  3.85 ,  2.32 , 18.61 ,  1.  ,  1.  ,  4.  ,  1.  ],
##        [ 6.  , 258.  , 110.  ,  3.08 ,  3.215, 19.44 ,  1.  ,  0.  ,  3.  ,  1.  ],
##        [ 8.  , 360.  , 175.  ,  3.15 ,  3.44 , 17.02 ,  0.  ,  0.  ,  3.  ,  2.  ],
##        [ 6.  , 225.  , 105.  ,  2.76 ,  3.46 , 20.22 ,  1.  ,  0.  ,  3.  ,  1.  ],
##        [ 8.  , 360.  , 245.  ,  3.21 ,  3.57 , 15.84 ,  0.  ,  0.  ,  3.  ,  4.  ],
##        [ 4.  , 146.7  ,  62.  ,  3.69 ,  3.19 , 20.   ,  1.  ,  0.  ,  4.  ,  2.  ],
##        [ 4.  , 140.8  ,  95.  ,  3.92 ,  3.15 , 22.9  ,  1.  ,  0.  ,  4.  ,  2.  ],
##        [ 6.  , 167.6  , 123.  ,  3.92 ,  3.44 , 18.3  ,  1.  ,  0.  ,  4.  ,  4.  ],
##        [ 6.  , 167.6  , 123.  ,  3.92 ,  3.44 , 18.9  ,  1.  ,  0.  ,  4.  ,  4.  ],
##        [ 8.  , 275.8  , 180.  ,  3.07 ,  4.07 , 17.4  ,  0.  ,  0.  ,  3.  ,  3.  ],
##        [ 8.  , 275.8  , 180.  ,  3.07 ,  3.73 , 17.6  ,  0.  ,  0.  ,  3.  ,  3.  ],
```