

Doing More with Data: numpy

Introduction to Python

Data Sciences Institute, University of Toronto

Instructor: A Mahfouz | TA: Kaylie Lau

July 2022

Contents:

1. Setup

2. numpy

A. Intro and arrays

B. Loading numpy

C. Creating arrays

D. Reshaping

E. Basic operations

F. Logic and filtering

G. Loading data

Setup

Anaconda

Install `numpy` if you haven't already. If you're not sure, go through the steps below.

1. Open Anaconda Prompt. You may have to do this in admin mode.
2. Type `conda install numpy` and press enter to run.
3. Press `enter` to answer `Y` to any prompts.

Google Colab

If you're using Colab, you're all set! `numpy` is already installed.

numpy

What is numpy ?

numpy is a Python package for scientific computing. It gives us *homogenous multidimensional arrays* -- a way of representing grids of elements where all elements are the same data type -- and functions to work efficiently with them. numpy both underpins and complements other data science libraries, including pandas, matplotlib, and scikit-learn.

Arrays vs lists

numpy arrays	Python list
all elements must be the same type	elements can be different types
fixed size	can change size
n-dimensional	1-dimensional
faster to process	slower to process
consumes less memory	consumes more memory

Loading numpy

We can `import` numpy like any other module. For convenience, numpy is typically loaded as `np` -- an alias that makes referencing it easier.

```
In [1]: import numpy as np
```

numpy arrays

The main object in `numpy` is the `ndarray`, also referred to as the `array`. Dimensions in an array are called *axes*. Most arrays we'll work with will be one-dimensional *vectors* and two-dimensional *matrices*.

We can create an array by calling `np.array()` and passing in data as a single value, like a list. Below is a matrix. The first axis has a length of two, and the second axis has a length of 3.

```
In [2]: a = np.array([[1, 2, 3],  
                      [3, 2, 1]])  
a
```

Out[2]:

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

An array has an `ndim` attribute indicating the number of its axes, a `size` indicating the number of values it has, and a `shape` indicating its size in each dimension. It also has a `dtype` describing what data type all of the elements in the array are.

```
In [3]: # notice that the shape is rows x columns
        print(a.shape)

#int32 is a numpy-provided dtype
print(a.dtype)
```

```
(2, 3)
int32
```

We can create arrays with placeholder content in several ways. This is useful when we know how many elements will be in an array, but not their values, as `numpy` arrays have fixed size. The full list is in [**numpy's documentation**](#).

```
In [4]: # create an 2x3x2 array of zeros. notice the double parentheses
        np.zeros((2, 3, 2))
```

Out[4]:

```
array([[[0., 0.],
        [0., 0.],
        [0., 0.]],
       [[0., 0.],
        [0., 0.],
        [0., 0.]])
```

```
In [5]: # create an array of ones based on the earlier a array
        np.ones_like(a)
```

Out[5]:

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

We can also create arrays by specifying a range of values through `arange()` or generating random ones through functions like `random.randint()` and `random.random()`.

```
In [6]: # create a 1D array from 1 til 10 in steps of 2
        np.arange(1, 10, 2)
```

Out[6]:

```
array([1, 3, 5, 7, 9])
```

```
In [7]: # seed is used to create a reproducible random example
        np.random.seed(1)
        # create a 3x4 array of random integers between 1 and 10
        np.random.randint(1, 10, (3, 4))
```

Out[7]:

```
array([[6, 9, 6, 1],
       [1, 2, 8, 7],
       [3, 5, 6, 3]])
```

We can repeat values and arrays to create bigger ones with `repeat()` and `tile()`.

```
In [8]: # create a 1D array through repetition
        np.repeat(10, 5)
```

Out[8]:

```
array([10, 10, 10, 10, 10])
```

```
In [9]: # create a 2D array through repetition
        onedim_arr = np.array([1, 2, 3, 4, 5])
        multidim_arr = np.tile(onedim_arr, (5,1))
        multidim_arr
```

Out[9]:

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

```
In [10]: # repeat onedim_arr in multiple directions
         another_arr = np.tile(onedim_arr, (3,3))
         another_arr
```

Out[10]:

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

Simulating data

Sometimes, it can be useful to simulate data to make sure an analysis works as expected.

numpy's `random` submodule provides support for drawing random samples from a variety of distributions. To generate random samples, we first call `default_rng()` to create a sample Generator. Then, we use the Generator's various distribution methods, like `normal()`, to create sample arrays.


```
In [11]: # create the sample Generator
         rng = np.random.default_rng(seed=42)

         # create a 5x6 array w/ normally distributed data
         # distribution has mean 10, standard dev. 2.5
         randrng = rng.normal(10, 2.5, (5, 4))
         randrng
```

Out[11]:

```
array([[10.7617927 ,  7.40003973, 11.87612799, 12.35141179],
       [ 5.12241203,  6.74455123, 10.31960101,  9.20939352],
       [ 9.95799711,  7.86739018, 12.19849494, 11.94447984],
       [10.16507674, 12.81810302, 11.16877336,  7.85176884],
       [10.92187696,  7.6027935 , 12.19612575,  9.87518522]])
```

```
In [12]: # using the same Generator to draw a random sample from a Poisson distribution
         poissonrng = rng.poisson(1, (5, 6))
         poissonrng
```

Out[12]:

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

Reshaping arrays

We can change the shape of an array in various ways, leaving the elements the same.

```
In [13]: # get the transpose of an array
         b = np.array([[5, 4, 3, 2, 1, 0],
                       [10, 8, 6, 4, 2, 0]])
         b.T
```

Out[13]:

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

```
In [14]: # change dimensions
         b.reshape(4, 3)
```

Out[14]:

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

If two arrays share the same size along an axis, we can stack them with `np.hstack()` and `np.vstack()`. Notice that we pass the arrays to stack as a tuple.

```
In [15]: # stack a and b horizontally
         np.hstack((b, a))
```

Out[15]:

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

```
In [16]: # reshape b before stacking vertically
         np.vstack((b.reshape(4, 3), a))
```

Out[16]:

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

Basic operations

numpy arrays let us perform vector operations, manipulating all the elements in an axis without writing loops.

```
In [17]: arr1 = np.array([5, 10, 15, 20])  
         arr1
```

```
Out[17]:  
array([ 5, 10, 15, 20])
```

```
In [18]: arr2 = np.arange(5, 9)  
         arr2
```

```
Out[18]:  
array([5, 6, 7, 8])
```

We can perform operations when arrays are the same length along the axis in use, or when values can be *broadcast*, or repeated, along an axis.

```
In [19]: # arr1 and arr2 are the same length  
        arr1 - arr2
```

Out[19]:

```
array([ 0,  4,  8, 12])
```

```
In [20]: # multiple each element in arr1 by 2  
        arr1 * 2
```

Out[20]:

```
array([10, 20, 30, 40])
```

```
In [21]: # incompatible shapes
         arr2 + np.array([1, 2])
```

```
-----
ValueError                                Traceback (most recent call last)
Input In [21], in <cell line: 2>()
      1 # incompatible shapes
----> 2 arr2 + np.array([1, 2])
```

ValueError: operands could not be broadcast together with shapes (4,) (2,)

We can also summarize the values in an array.

```
In [22]: print(f'''arr1 sums to {arr1.sum()}.  
            Its max value is {arr1.max()}, and its mean is {arr1.mean()}.''' )
```

arr1 sums to 50.

Its max value is 20, and its mean is 12.5.

Some descriptive statistics are only numpy functions, and are not available as array methods like we saw above.

```
In [23]: np.median(arr1)
```

```
Out[23]:
```

```
12.5
```

```
In [24]: arr1.median()
```

AttributeError

Traceback (most recent call last)

Input **In [24]**, in <cell line: 1>()

----> 1 arr1.median()

AttributeError: 'numpy.ndarray' object has no attribute 'median'

Operations in multiple dimensions

In a multi-dimensional array like a matrix, elements of same-sized arrays will be paired up for operations, just as with vectors. If we perform an operation with a matrix and a vector of the same size in one dimension, the vector can be *broadcast* to repeat along the other dimension.

```
In [25]: tens = np.arange(0, 120, 10).reshape(3, 4)
         tens
```

Out[25]:

```
array([[ 0, 10, 20, 30],
       [ 40, 50, 60, 70],
       [ 80, 90, 100, 110]])
```

```
In [26]: horizontal = np.array([-5, -10, -15, -20])  
  
tens + horizontal
```

Out[26]:

```
array([[ -5,  0,  5, 10],  
       [35, 40, 45, 50],  
       [75, 80, 85, 90]])
```

```
In [27]: vertical = np.array([[100],  
                               [200],  
                               [300]])  
  
tens + vertical
```

Out[27]:

```
array([[100, 110, 120, 130],  
       [240, 250, 260, 270],  
       [380, 390, 400, 410]])
```

We can still calculate statistics for multidimensional arrays, but we must specify the axis to calculate over. To calculate values for each column, we use `axis=0`. To calculate for each row, we use `axis=1`.

```
In [28]: tens.mean(axis=0)
```

Out[28]:

```
array([40., 50., 60., 70.])
```

```
In [29]: tens.mean(axis=1)
```

Out[29]:

```
array([15., 55., 95.])
```

Indexing, slicing, and iterating

We can index and slice arrays like we would a list.

```
In [30]: arr1
```

```
Out[30]:
```

```
array([ 5, 10, 15, 20])
```

```
In [31]: arr1[1]
```

```
Out[31]:
```

```
10
```

```
In [32]: arr1[1:3]
```

```
Out[32]:
```

```
array([10, 15])
```

We can iterate over arrays as well, though vectorized numpy operations are preferred when possible.

```
In [33]: for i in arr1:  
         print(i)
```

```
5  
10  
15  
20
```

Multidimensional arrays like matrices have one index per axis. We can pass in more than one index within the square brackets.

```
In [34]: tens
```

```
Out[34]:
```

```
array([[ 0, 10, 20, 30],  
       [ 40, 50, 60, 70],  
       [ 80, 90, 100, 110]])
```

```
In [35]: # indexing goes row, column  
tens[1, 2]
```

```
Out[35]:
```

```
60
```

```
In [36]: # get the first row  
tens[0]
```

```
Out[36]:
```

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

```
In [37]: # get the first column  
tens[:,0]
```

```
Out[37]:
```

```
array([ 0, 40, 80])
```

```
In [38]: # slice rows 1-2, columns 2-3  
tens[0:2, 1:3]
```

```
Out[38]:
```

```
array([[10, 20],  
       [50, 60]])
```


Mutations and copies

We can also update individual elements in an array. Note that any variables that refer to that array will also change, just like with mutating lists. To make an independent copy of an array, use the `.copy()` method.

```
In [39]: # create a 3x4 array of random integers
        matrix = np.random.randint(1, 11, 12).reshape(3, 4)
matrix2 = matrix

#make a copy
matrix3 = matrix.copy()

matrix
```

Out[39]:

```
array([[ 5,  3,  5,  8],
       [ 8, 10,  2,  8],
       [ 1,  7, 10, 10]])
```

```
In [40]: # replace the second row
        matrix2[1] = [0, 0, 0, 0]
matrix2
```

Out[40]:

```
array([[ 5,  3,  5,  8],
       [ 0,  0,  0,  0],
       [ 1,  7, 10, 10]])
```

```
In [41]: # the original also changed
        matrix
```

Out[41]:

```
array([[ 5,  3,  5,  8],
       [ 0,  0,  0,  0],
       [ 1,  7, 10, 10]])
```

```
In [42]: # the copy did not
        matrix3
```

Out[42]:

```
array([[ 5,  3,  5,  8],
       [ 8, 10,  2,  8],
       [ 1,  7, 10, 10]])
```

To filter use a boolean expression as a mask, pass it into square brackets after the array to mask.

Logic and filtering

numpy arrays work with boolean expressions. Each element is checked, and the result is an array of True / False values. They resulting arrays sometimes called *masks* because they are used to mask, or filter, data.

```
In [43]: tens
```

```
Out[43]:
```

```
array([[ 0, 10, 20, 30],
       [ 40, 50, 60, 70],
       [ 80, 90, 100, 110]])
```

```
In [44]: # evaluate whether each element is divisible by 3
        tens % 3 == 0
```

```
Out[44]:
```

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

```
In [45]: # the same thing with for loops
        masked = []

        for row in tens:
            masked_row = []
            for col in row:
                masked_row.append(col % 3 == 0)
            masked.append(masked_row)

masked
```

Out[45]:

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

```
In [46]: tens[tens % 3 == 0]
```

Out[46]:

```
array([ 0, 30, 60, 90])
```

```
In [47]: # also works
        mask = tens % 3 == 0
        tens[mask]
```

Out[47]:

```
array([ 0, 30, 60, 90])
```

```
In [48]: # comparison in standard python
         filtered_data = []

         for row in tens:
             for col in row:
                 if col % 3 == 0:
                     filtered_data.append(col)

         filtered_data
```

Out[48]:

```
[0, 30, 60, 90]
```


We can even use masks to generate new arrays with conditionals. `np.where()` takes as its arguments a boolean expression, an expression to evaluate if `True`, and an expression to evaluate otherwise. This is analagous to creating a new array based on an old one with a `for` loop and `if/else` statements, but much faster.

```
In [49]: np.where(tens % 3 == 0, # condition
                  tens, # return the element if True
                  0) # return 0 if False
```

Out[49]:

```
array([[ 0,  0,  0, 30],
       [ 0,  0, 60,  0],
       [ 0, 90,  0,  0]])
```

```
In [50]: result = []

for row in tens:
    result_row = []
    for col in row:
        if col % 3 == 0:
            result_row.append(col)
        else:
            result_row.append(0)
    result.append(result_row)

result
```

Out[50]:

```
[[0, 0, 0, 30], [0, 0, 60, 0], [0, 90, 0, 0]]
```

Loading flat files to numpy arrays

We can also load data from files into numpy arrays. Recall the California housing csv we loaded earlier:

```
In [51]: with open('sample_data/california_housing_test.csv', 'r') as f:
          # print the first five lines
          for i in range(5):
              print(f.readline())
```

```
longitude,latitude,housing_median_age,total_rooms,total_bedrooms,population,households,median_income,median_house_value
```

```
-122.05,37.37,27,3885,661,1537,606,6.6085,344700
```

```
-118.3,34.26,43,1510,310,809,277,3.599,176500
```

```
-117.81,33.78,27,3589,507,1484,495,5.7934,270500
```

```
-118.36,33.82,28,67,15,49,11,6.1359,330000
```

We can load this data to a numpy array using `genfromtxt()`, which takes a path to a file as an argument. Optional additional arguments include the `delimiter`, indicating how values are separated; `names`, which we can use to tell numpy that the first row contains column names; `skip_header`, which we can use to skip the first few lines; and arguments for how to handle missing values.

```
In [56]: housing_data = np.genfromtxt('sample_data/california_housing_test.csv',
                                     delimiter=',',
                                     #skip_header=1, # Let's use names instead
                                     names=True
                                     )
housing_data
```

Out[56]:

```
array([(-122.05, 37.37, 27., 3885., 661., 1537., 606., 6.6085, 344700.),
       (-118.3 , 34.26, 43., 1510., 310., 809., 277., 3.599 , 176500.),
       (-117.81, 33.78, 27., 3589., 507., 1484., 495., 5.7934, 270500.),
       ...,
       (-119.7 , 36.3 , 10., 956., 201., 693., 220., 2.2895, 62000.),
       (-117.12, 34.1 , 40., 96., 14., 46., 14., 3.2708, 162500.),
       (-119.63, 34.42, 42., 1765., 263., 753., 260., 8.5608, 500001.)],
      dtype=[('longitude', '<f8'), ('latitude', '<f8'), ('housing_median_age', '<f8'), ('total_rooms', '<f8'), ('total_bedrooms', '<f8'), ('population', '<f8'), ('households', '<f8'), ('median_income', '<f8'), ('median_house_value', '<f8')])
```

If you look closely, the array we got back is not a matrix. Each row looks like a tuple with comma-separated values. If we check the `shape`, we see 3000 rows and what looks like no columns. The `dtype` attribute lists all our field names.

```
In [57]: housing_data.shape
```

```
Out[57]:
```

```
(3000,)
```

```
In [58]: housing_data.dtype
```

```
Out[58]:
```

```
dtype([('longitude', '<f8'), ('latitude', '<f8'), ('housing_median_age', '<f8'), ('total_rooms', '<f8'), ('total_bedrooms', '<f8'), ('population', '<f8'), ('households', '<f8'), ('median_income', '<f8'), ('median_house_value', '<f8')])
```

In this case, `genfromtxt` returned a *structured array*, a different type of array than the ones we have seen so far. We can refer to fields by putting the field name as a string in square brackets, similar to referencing a dictionary key. However, we will soon see a data type in another package, `pandas`, that is even better suited to working with columns in tabular data like this.

```
In [59]: np.median(housing_data['housing_median_age'])
```

```
Out[59]:
```

```
29.0
```

References

- NumPy. *Basic Numpy*. [**https://numpy.org/devdocs/user/quickstart.html**](https://numpy.org/devdocs/user/quickstart.html)
- NumPy. *Numpy Routines*. [**https://numpy.org/doc/stable/reference/routines.html**](https://numpy.org/doc/stable/reference/routines.html)