

Intro. Comp. for Data Science (FMI08)

Dr. Nono Saha

April 26, 2023

Max Planck Institute for Mathematics in the Sciences
University of Leipzig/ScaDS.AI

Spring 2023

1. NumPy numerics
2. NumPy - Advanced indexing
3. NumPy - Broadcasting
4. NumPy - Basic file I/O
5. Structure of a Data Science (ML) project
6. Homework 3

NumPy numerics

NumPy numerics: basic operations

All basic mathematical operators in **Python** are implemented for arrays. They are applied element-wise to the array values.

```
1 np.arange(3) + np.  
  arange(3)  
2     ## array([0, 2, 4])  
3  
4 np.arange(3) - np.  
  arange(3)  
5     ## array([0, 0, 0])  
6  
7 np.arange(3) + 2  
8     ## array([2, 3, 4])  
9
```

```
np.arange(3) * np.arange  
(3)  
     ## array([0, 1, 4])  
  
np.arange(1,4)/np.arange  
(1,4)  
     ## array([1., 1., 1.])  
  
np.arange(3) * 3  
     ## array([0, 3, 6])
```

```
1 np.full((2,2), 2) ** np.arange(4).reshape((2,2))  
2 np.full((2,2), 2) ** np.arange(4)  
3     ## Which of the two instructions will work?  
4
```

NumPy numerics: mathematical functions

The package provides a wide variety of basic mathematical functions that are vectorized. In general, they will be faster than their base equivalents (e.g. `np.sum()` vs `sum()`).

```
1 np.sum(np.arange(1000))
```

```
2 ## 499500
```

```
3  
4 np.cumsum(np.arange(10))
```

```
5 ## array([ 0,  1,  3,  6, 10, 15, 21, 28, 36, 45])
```

```
6  
7 np.log10(np.arange(1,11))
```

```
8 ## array([0., 0.30103, 0.47712125, 0.60205999, 0.69897,
```

```
9 ## 0.77815125, 0.84509804, 0.90308999, 0.95424251, 1. ])
```

```
10  
11 np.median(np.arange(10))
```

```
12 ## 4.5
```

NumPy numerics: matrix multiplication

It is supported using the `matmul()` function or the `@` operator,

```
1  x = np.arange(6).reshape(3,2)
2  y = np.tri(2,2)
3  x @ y
4  ## array([[1., 1.], [5., 3.], [9., 5.]])
5
6  y.T @ y
7  ## array([[2., 1.], [1., 1.]])
8
9  np.matmul(x.T, x)
10 ## array([[20, 26], [26, 35]])
11
12 y @ x
13 ## Can this work?
14
```

NumPy numerics: other linear algebra functions

The standard linear algebra functions are (mostly) implemented in the `linalg` submodule. See [here](#) for more details.

```
1  np.linalg.det(y)
2  ## 1.0
3
4  np.linalg.eig(x.T @ x)
5  ## (array([ 0.43988174, 54.56011826]), array([[ -0.79911221,
6  -0.6011819 ],
7  ## [ 0.6011819 , -0.79911221]]))
8
9  np.linalg.inv(x.T @ x)
10 ## array([[ 1.45833333, -1.08333333], [-1.08333333,
11 0.83333333]])
12
13 np.linalg.cholesky(x.T @ x)
14 ## array([[4.47213595, 0.],[5.81377674, 1.09544512]])
```

NumPy numerics: random values

NumPy has another submodule called random for functions used to generate random values,

To use this, you should construct a generator via `default_rng()`, with or without a seed, and then use the generator's methods to obtain your desired random values.

```
1  rng = np.random.default_rng(seed = 1234)
2  rng.random(3) # ~ Uniform [0,1)
3  ## array([0.97669977, 0.38019574, 0.92324623])
4
5  rng.normal(0, 2, size = (2,2))
6  ## array([[ 0.30523839,  1.72748778],[ 5.82619845,
7  -2.95764672]])
8
9  rng.binomial(n=5, p=0.5, size = 10)
10 ## array([2, 4, 2, 2, 3, 4, 4, 3, 3, 3])
```


NumPy - Advanced indexing

From last time: subsetting with tuples

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

```
1  x = np.arange(6)
2  x
3  ## array([0, 1, 2, 3, 4, 5])
4
5  x[(0,1,3),]
6
7  ## array([0, 1, 3])
8
9  x[(0,1,3)]
10
11 ## Traceback (most recent call last):
12 File "<stdin>", line 1, in <module>
13 IndexError: too many indices for array: array is 1-
14 dimensional, but three were indexed
```

Question

What if we use the list instead?

NumPy - Advanced indexing: exercise

Given the following matrix,

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

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

NumPy - Advanced indexing: boolean indexing

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

```
1 x = np.arange(6)
2 ## array([0, 1, 2, 3, 4, 5])
3
4 x[[True, False, True, False, True, False]]
5 ## array([0, 2, 4])
6
7 x[np.array([True, True, False, False, True, False])]
8 ## array([0, 1, 4])
9
```

The utility comes from vectorized comparison operations,

```
1 x > 3
2 ## array([False, False, False, False,  True,  True])
3 x[x>3]
4 ## array([4, 5])
5 x % 2 == 1
6 ## array([False,  True, False,  True, False,  True])
7
```

NumPy - Advanced indexing: 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.

```
1  x = np.arange(6)
2  x
3  ## array([0, 1, 2, 3, 4, 5])
4
5  y = x % 2 == 0
6  y
7  ## array([ True, False,  True, False,  True, False])
8
9  ~y
10 ## array([False,  True, False,  True, False,  True])
11
12 y & (x > 3)
13 ## array([False, False, False, False,  True, False])
14
15 y | (x > 3)
16 ## array([ True, False,  True, False,  True,  True])
17
```

NumPy - meshgrid

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

```
1 pts = np.arange(3)
2 x, y = np.meshgrid(pts, pts)
3 x
4 ## array([[0, 1, 2], [0, 1, 2], [0, 1, 2]])
5
6 y
7 ## array([[0, 0, 0], [1, 1, 1], [2, 2, 2]])
8
9 np.sqrt(x**2 + y**2)
10
11 ## array([[0.          , 1.          , 2.          ],
12 ##        [1.          , 1.41421356, 2.23606798],
13 ##        [2.          , 2.23606798, 2.82842712]])
14
```

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 of 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

NumPy - Broadcasting: 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

- they are equal, or
- 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.

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

```
      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])
```

NumPy - Broadcasting: mechanism

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

```
2 y = 1
```

```
3 x+y
```

```
4  
5 x      (2d array): 4 x 3
```

```
6 y      (1d array):      1
```

```
7 -----  
8 x+y    (2d array): 4 x 3
```

```
9  
10  
11 x = np.arange(12).reshape  
   ((4,3))
```

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

```
13 x+y
```

```
14  
15 x      (2d array): 4 x 3
```

```
16 y      (1d array):      3
```

```
17 -----  
18 x+y    (2d array): 4 x 3
```

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

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

```
x+y
```

```
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
```

```
x      (2d array): 3 x 4
```

```
y      (1d array): 3 x 1
```

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

NumPy - Broadcasting: example for data 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()`.

```
1  rng = np.random.default_rng(1234)
2  d = rng.normal(loc=[-1,0,1], scale=[1,2,3], size=(1000,3))
3  d.mean(axis=0)
4  ## array([-1.0294382 , -0.01396257,  1.01241784])
5
6  d.std(axis=0)
7  ## array([0.99674719,  2.03222595,  3.10625219])
8
```

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

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

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

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

NumPy - Basic file I/O

NumPy - Basic file I/O: reading and writing arrays

We will not spend much time on this as most data you will encounter is more likely to be in 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.

```
1 x = np.arange(1e5)
2 np.save("data/x.npy", x)
3 new_x = np.load("data/x.npy")
4 np.all(x == new_x)
5
6 ## True
7
```

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

If you need to read delimited (CSV, tsv, etc.) data into a **NumPy** array, you can use `genfromtxt()`.