

Tools and Programming Languages for Data Science

NumPy

Study Program Data Science
Prof. Dr. Tillmann Schwörer

Python foundations

Data types
Operators
Functions
Control flow and iterators
Programming concepts & paradigms



See also Precourse
Programming

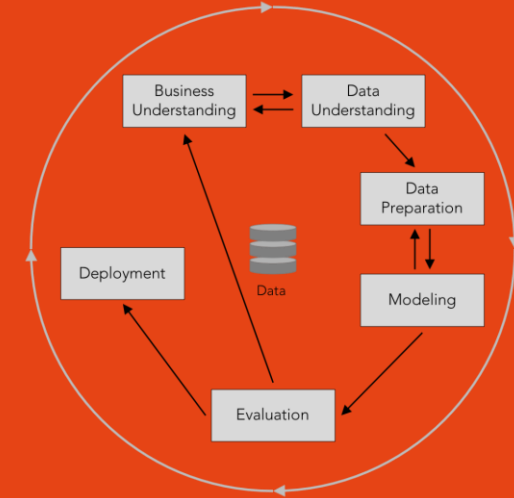
Tooling

Installation
Visual Studio Code
Jupyter Notebooks
Packages
Virtual Environments
Git and Github



Python

Data Science Workflow



NumPy



pandas

matplotlib

Our course agenda

- ▶ **Introduction and overview**
- ▶ **NumPy**: Basic data handling with Numpy arrays
- ▶ **Pandas**
 - ◆ Exploratory data analysis
 - ◆ Data consolidation
 - ◆ Data cleaning
- ▶ **Data visualization using Matplotlib and Seaborn**
- ▶ **Interacting with APIs**
- ▶ **Interacting with SQL databases**
- ▶ **Version Control with Git and GitHub**
- ▶ **Advanced Python**

Characteristics of NumPy arrays

► **Multidimensional:**

- ◆ Vectors (1 dimension), matrices (2), tensors (3 or more)
- ◆ We say „dimensions“ or „axes“

► **Homogeneous data type**

- ◆ All elements must be of the same data type: e.g. only integers
- ◆ Enables efficient storage and fast computations because there is no need for type checking

► **Array attributes**

- ◆ ndim: number of dimensions (e.g. 1, 2, 3, ...)
- ◆ shape: length along each dimension (e.g. (50, 5) for data with 50 rows and 5 cols)
- ◆ size: total number of elements (e.g. 250 for data with 50 rows and 5 cols)
- ◆ dtype: data type (e.g. int32, float64, ...)

NumPy Key Features

- ▶ **Multidimensional arrays**
 - ◆ 1, 2, or higher-dimensional data
- ▶ **Efficient indexing, slicing and manipulation of array elements**
- ▶ **Mathematical/statistical functions** operating on arrays:
 - ◆ mean
 - ◆ max
 - ◆ log
 - ◆ ...
- ▶ **Linear Algebra**
 - ◆ dot product
 - ◆ solving system of equations
 - ◆ ...

N-dimensional array

1D array

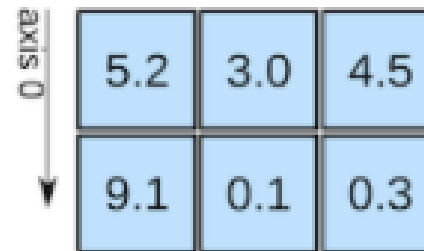


axis 0 →

shape: (4,)

```
np.array([7, 2, 9, 10])
```

2D array



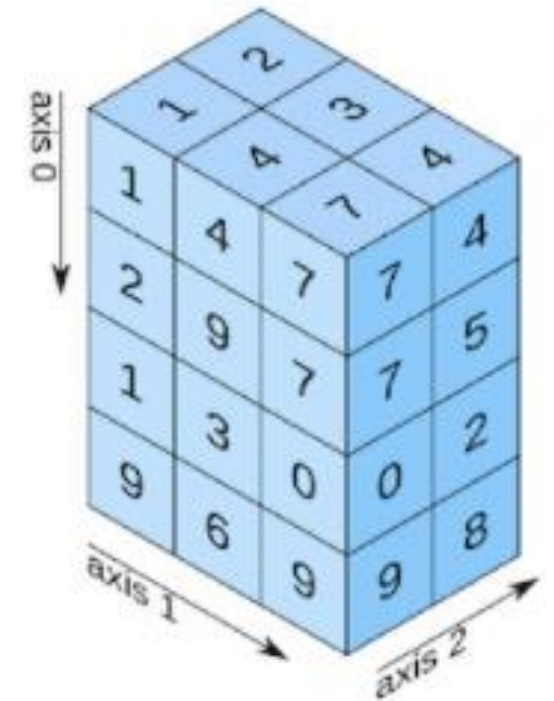
axis 0 ↓

axis 1 →

shape: (2, 3)

```
np.array([(5.2, 3.0, 4.5),  
          ... (9.1, 0.1, 0.3)])
```

3D array



axis 0 ↓

axis 1 →

axis 2 →

shape: (4, 3, 2)

Advantages of NumPy

- ▶ **Highly performant**
- ▶ **Memory efficient**
- ▶ **Broad range of operations relevant for data analysis and manipulation**
- ▶ **Linked to or interoperable with many other Python packages**
 - ◆ Pandas Series and DataFrames are built on numpy
 - ◆ Machine learning with scikit-learn expects numpy arrays as input,
 - ◆ Arrays can be visualized with matplotlib
 - ◆ Pytorch tensors are interoperable with numpy arrays
 - ◆ ...

Characteristics of NumPy arrays

▶ „Size-immutable“ / „fixed-size“:

- ◆ Once the numpy array is created its array size cannot be changed → the number of elements cannot be changed
- ◆ Append, insert, remove leads to creation of a new array

▶ “Element-mutable”

- ◆ Elements within the array can be modified in place → no new array is created
- ◆ Modification of individual elements, slices, or more complex selections (e.g. based on some boolean criterium)

List vs NumPy array

Lists

- Built-in
- 1-dimensional
- Heterogeneous items
- Flexible resizing: adding or deleting elements can often occur in place
- Not optimized for Math

Numpy Arrays

- Not built-in
- Multi-dimensional
- Homogeneous items
- Fixed size: adding or deleting items creates a new array
- High performance math

Why are Numpy arrays (usually) faster than lists?

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

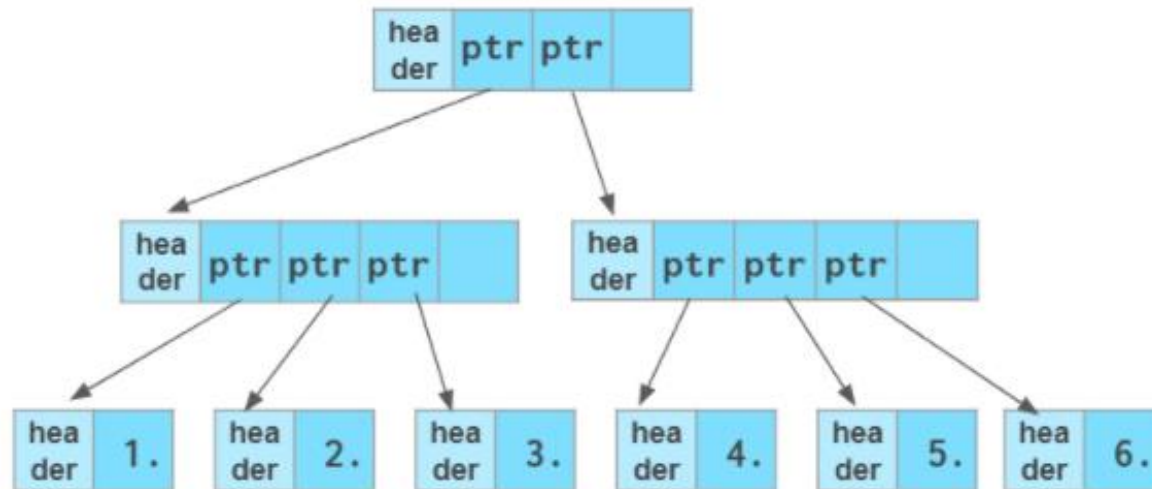
python list

1.	2.	3.
4.	5.	6.

vs

numpy array

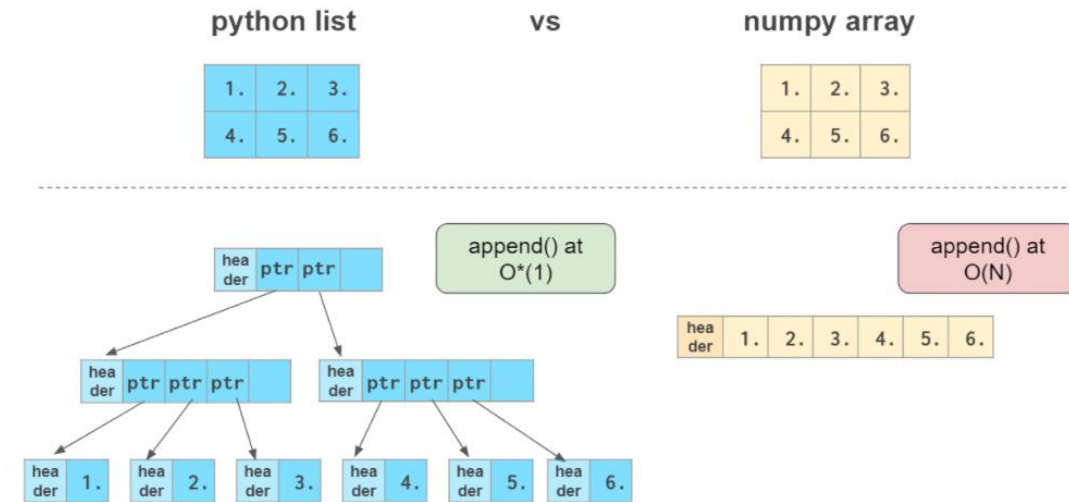
1.	2.	3.
4.	5.	6.



header	1.	2.	3.	4.	5.	6.
--------	----	----	----	----	----	----

Why are Numpy arrays (usually) faster than lists?

- ▶ Use **contiguous memory**
- ▶ Each element needs **less memory**
- ▶ **No type checking** needed when iterating through object
- ▶ **Vectorized operations:** For typical operations on arrays, we can avoid Python loops (slow) → Numpy implicitly delegates execution of loops to underlying C code



But: Slower if the size of the object changes (insert, append, delete)

Element-wise operations

```
squared = []  
for number in [2, 3]:  
    squared.append(number**2)
```

```
np.array([2, 3])**2
```

Vectorized operation: no
Python loop needed

a^2

=

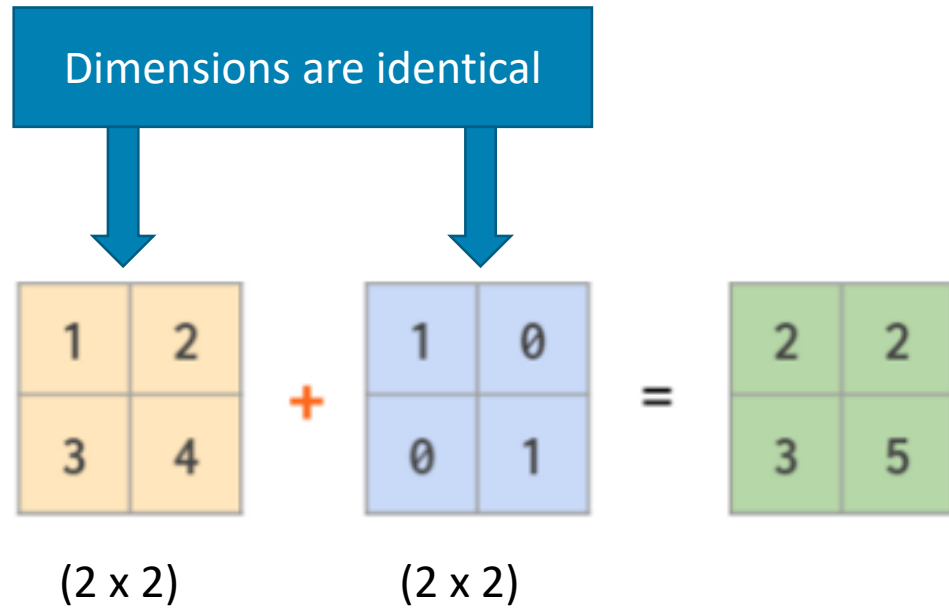
2	3
---	---

**** 2**

=

4	9
---	---

Element-wise operations



```
x = np.array([[1, 2], [3, 4]])  
y = np.array([[1, 0], [0, 1]])  
x + y
```


Broadcasting

Dimensions of the smaller array
are **broadcasted** to fit the
dimensions of the larger array

$$\begin{bmatrix} 1 & 2 \end{bmatrix} + 3 = \begin{bmatrix} 4 & 5 \end{bmatrix}$$

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

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} / \begin{bmatrix} 9 & 9 & 9 \\ 9 & 9 & 9 \\ 9 & 9 & 9 \end{bmatrix} = \begin{bmatrix} .1 & .2 & .3 \\ .4 & .5 & .7 \\ .8 & .9 & 1. \end{bmatrix}$$

Statistics and the axis Argument



Other statistics / aggregation functions

- min
- max
- mean
- median
- sum
- std
- var
- ...

Two alternatives

`np.max(`  `) =` 

The diagram shows a 1D array with elements 1, 2, and 3. The element 3 is highlighted in dark blue, indicating it is the maximum value. The array is labeled 'a' above it.

Top-level function of the numpy module

 `.max()` `=` 

The diagram shows a 1D array with elements 1, 2, and 3. The element 3 is highlighted in dark blue, indicating it is the maximum value.

Method of numpy arrays

Indexing and Slicing

`a = np.arange(1, 6)`

1	2	3	4	5
0	1	2	3	4

`a[2:4] = 0`



`a`

1	2	0	0	5
---	---	---	---	---

`a`

1	2	3	4
5	6	7	8
9	10	11	12

`a[1,2]`

1	2	3	4
5	6	7	8
9	10	11	12

`a[:,1:2]`

1	2	3	4
5	6	7	8
9	10	11	12

Similar to lists:

- ▶ Same indexing rules
- ▶ Applicable to any dimension
- ▶ Mutability: elements of an array can be changed in-place → potential side effects between multiple names pointing to the same numpy array

Boolean masking

a	1	2	3	4	5	6	7	6	5	4	3	2	1
----------	---	---	---	---	---	---	---	---	---	---	---	---	---

a > 5	False	False	False	False	False	True	True	True	False	False	False	False	False
-----------------	-------	-------	-------	-------	-------	------	------	------	-------	-------	-------	-------	-------


a[a > 5]	6	7	6
--------------------	---	---	---

Linear algebra

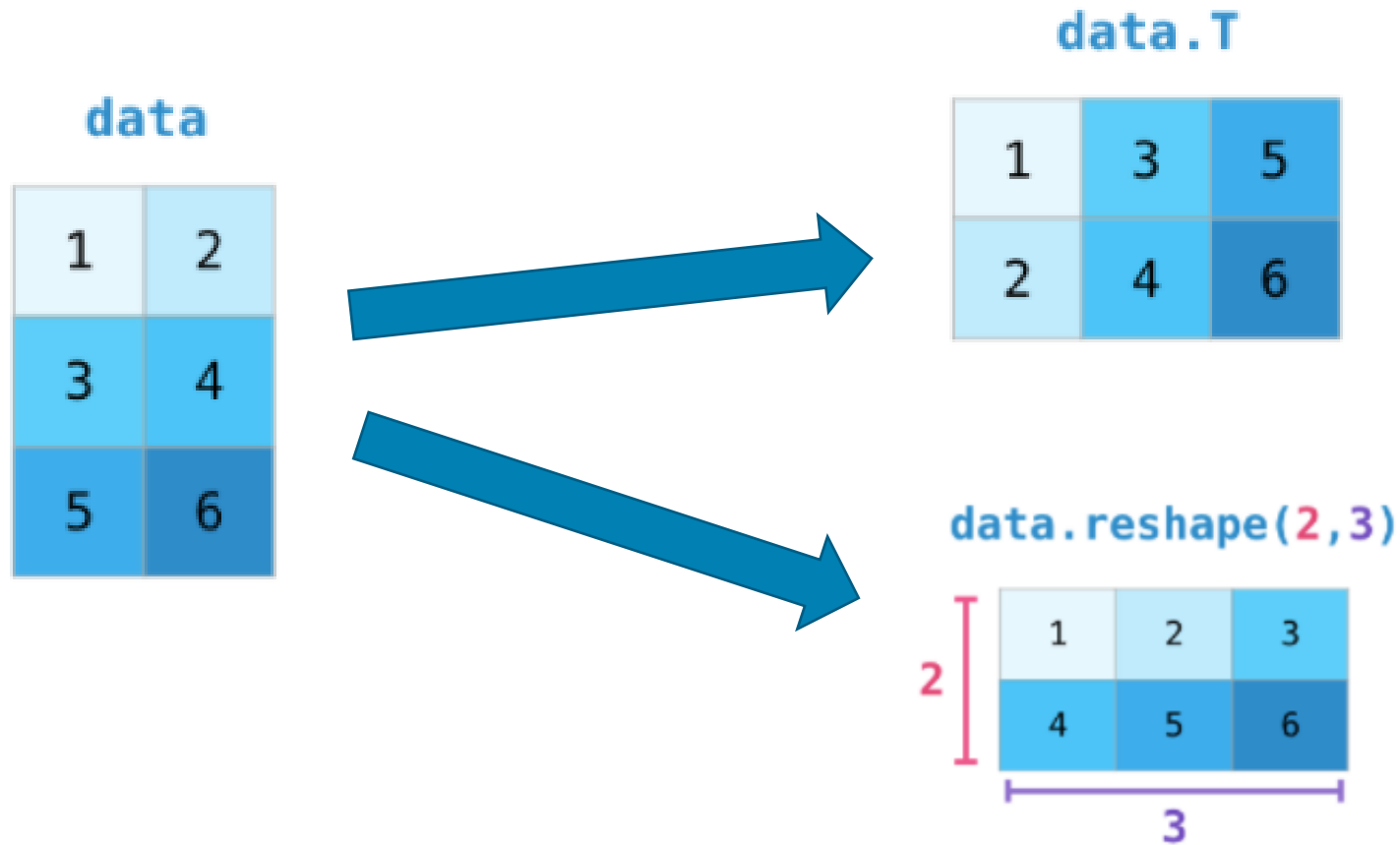
```
x = np.array([[1,2],[5,6],[7,8]])  
y = np.array([3,4])  
  
x.dot(y)  
x @ y
```

$$\begin{bmatrix} 1 & 2 \\ 5 & 6 \\ 7 & 8 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 4 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 \\ 5 & 6 \\ 7 & 8 \end{bmatrix}^{\begin{bmatrix} 3 & 4 \end{bmatrix}} = \begin{bmatrix} 1 \cdot 3 + 2 \cdot 4 \\ 5 \cdot 3 + 6 \cdot 4 \\ 7 \cdot 3 + 8 \cdot 4 \end{bmatrix} = \begin{bmatrix} 11 \\ 39 \\ 53 \end{bmatrix}$$

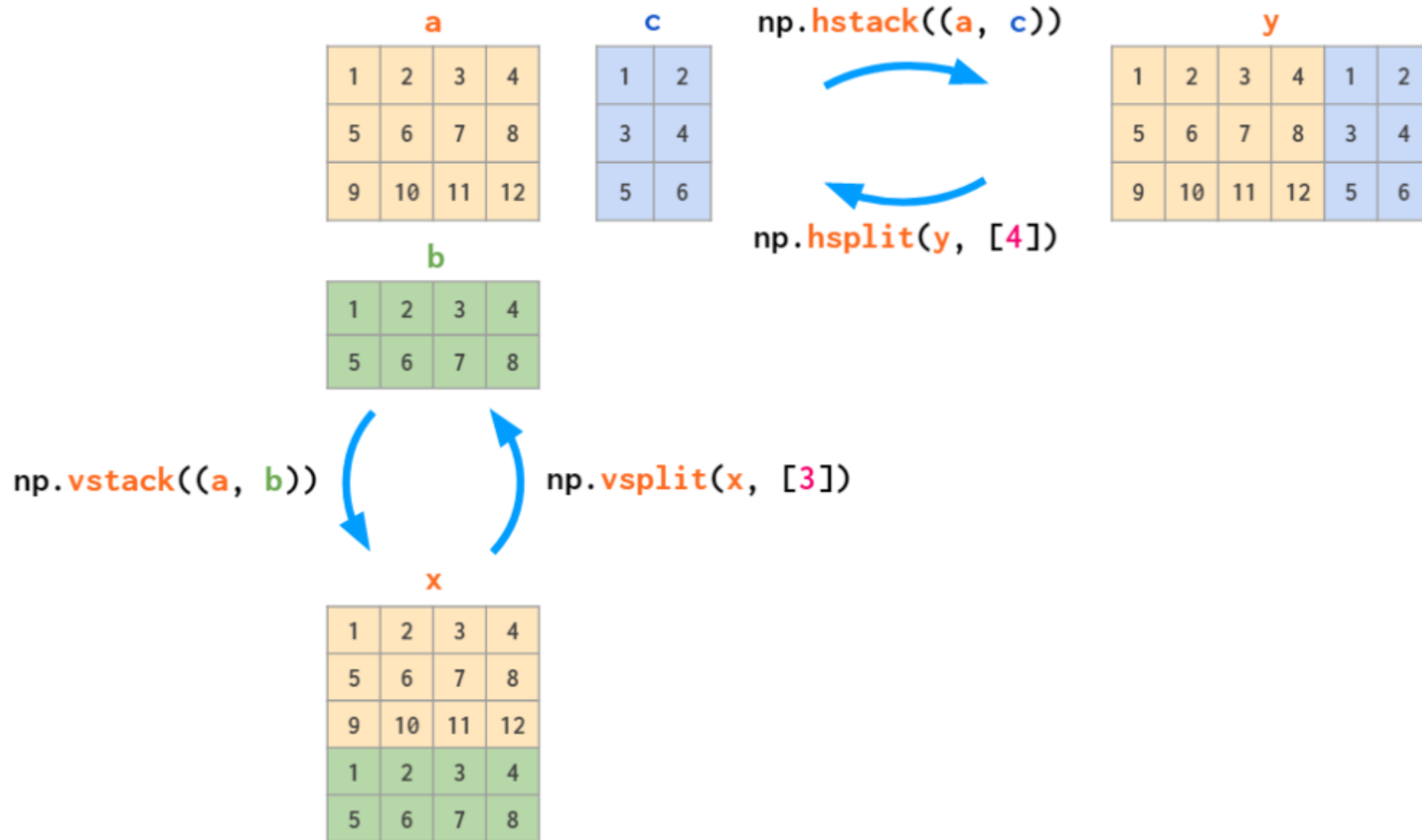
$(3 \times 2) \quad (2 \times 1)$ (3×1)



Transpose and reshape



Stacking and splitting



References

Visual guides to Numpy:

- <https://betterprogramming.pub/numpy-illustrated-the-visual-guide-to-numpy-3b1d4976de1d>
- <http://jalammar.github.io/visual-numpy/>