# SEBASTIAN RASCHKA



# Introduction to Artificial Neural Networks and Deep Learning

with Applications in Python

## Introduction to Artificial Neural Networks

with Applications in Python

Sebastian Raschka

#### **DRAFT**

Last updated: May 30, 2018

This book will be available at http://leanpub.com/ann-and-deeplearning.

Please visit https://github.com/rasbt/deep-learning-book for more information, supporting material, and code examples.

© 2016-2018 Sebastian Raschka

# **Contents**

F	Introduction to NumPy		
	F.1	N-dimensional Arrays	5
	F.2	Array Construction Routines	8
	F.3	Array Indexing	10
	F.4	Array Math and Universal Functions	12
	F.5	Broadcasting	15
	F.6	Advanced Indexing – Memory Views and Copies	17
	F.7	Comparison Operators and Mask	20
	F.8	Random Number Generators	22
	F.9	Reshaping Arrays	25
	F.10	Linear Algebra with NumPy Arrays	27
	F.11	Set Operations	32
	F.12	Serializing NumPy Arrays	34
	F.13	Conclusion	35

## Website

Please visit the GitHub repository to download the code examples accompanying this book and other supplementary material.

If you like the content, please consider supporting the work by buying a copy of the book on Leanpub. Also, I would appreciate hearing your opinion and feedback about the book, and if you have any questions about the contents, please don't hesitate to get in touch with me via mail@sebastianraschka.com. Happy learning!

Sebastian Raschka

## **About the Author**

Sebastian Raschka received his doctorate from Michigan State University developing novel computational methods in the field of computational biology. In summer 2018, he joined the University of Wisconsin–Madison as Assistant Professor of Statistics. Among others, his research activities include the development of new deep learning architectures to solve problems in the field of biometrics. Among his other works is his book "Python Machine Learning," a bestselling title at Packt and on Amazon.com, which received the ACM Best of Computing award in 2016 and was translated into many different languages, including German, Korean, Italian, traditional Chinese, simplified Chinese, Russian, Polish, and Japanese.

Sebastian is also an avid open-source contributor and likes to contribute to the scientific Python ecosystem in his free-time. If you like to find more about what Sebastian is currently up to or like to get in touch, you can find his personal website at https://sebastianraschka.com.

# Acknowledgements

I would like to give my special thanks to the readers, who provided feedback, caught various typos and errors, and offered suggestions for clarifying my writing.

- Appendix A: Artem Sobolev, Ryan Sun
- Appendix B: Brett Miller, Ryan Sun
- Appendix D: Marcel Blattner, Ignacio Campabadal, Ryan Sun, Denis Parra Santander
- Appendix F: Guillermo Monecchi, Ged Ridgway, Ryan Sun, Patric Hindenberger
- Appendix H: Brett Miller, Ryan Sun, Nicolas Palopoli, Kevin Zakka



## Appendix F

# Introduction to NumPy

This appendix offers a quick tour of the NumPy¹ library for working with multi-dimensional arrays in Python [Oliphant, 2006, Walt et al., 2011]. NumPy (short for Numerical Python) was created by Travis Oliphant in 2005, by merging Numarray into Numeric. Since then, the open source NumPy library has evolved into an essential library for scientific computing in Python and has become a building block of many other scientific libraries, such as SciPy, Scikit-learn, Pandas, and others. What makes NumPy so particularly attractive to the scientific community is that it provides a convenient Python interface for working with multi-dimensional array data structures efficiently; the NumPy array data structure is also called ndarray, which is short for *n*-dimensional array.

In addition to being mostly implemented in C and using Python as "glue language," the main reason why NumPy is so efficient for numerical computations is that NumPy arrays use contiguous blocks of memory that can be efficiently cached by the CPU. In contrast, Python lists are arrays of pointers to objects in random locations in memory, which cannot be easily cached and come with a more expensive memory-look-up. However, the computational efficiency and low-memory footprint come at a cost: NumPy arrays have a fixed size and are homogenous, which means that all elements must have the same type. Homogenous ndarray objects have the advantage that NumPy can carry out operations using efficient C loops and avoid expensive type checks and other overheads of the Python API. While adding and removing elements from the end of a Python list is very efficient, altering the size of a NumPy array is very expensive since it requires creating a new array and carrying over the contents of the old

<sup>&</sup>lt;sup>1</sup>http://www.numpy.org



array that we want to expand or shrink.

Besides being more efficient for numerical computations than native Python code, NumPy can also be more elegant and readable due to vectorized operations and broadcasting, which are features that we will explore in this appendix. While this appendix should be sufficient for following the code examples in this book if you are new to NumPy, there are many advanced NumPy topics that are beyond the scope of this book. If you are interested in a more in-depth coverage of NumPy, I selected a few resources that could be useful to you:

- Rougier, N.P., 2016. From Python to Numpy<sup>2</sup>
- Oliphant, T.E., 2015. A Guide to NumPy: 2nd Edition<sup>3</sup>. USA: Travis Oliphant, independent publishing.
- Varoquaux, G., Gouillart, E., Vahtras, O., Haenel, V., Rougier, N.P., Gommers, R., Pedregosa, F., Jędrzejewski-Szmek, Z., Virtanen, P., Combelles, C. and Pinte, D., 2015. SciPy Lecture Notes<sup>4</sup>
- The official NumPy documentation<sup>5</sup>

#### F.1 N-dimensional Arrays

NumPy is built around ndarrays<sup>6</sup> objects, which are high-performance multi-dimensional array data structures. Intuitively, we can think of a one-dimensional NumPy array as a data structure to represent a vector of elements – you may think of it as a fixed-size Python list where all elements share the same type. Similarly, we can think of a two-dimensional array as a data structure to represent a matrix or a Python list of lists. While NumPy arrays can have up to 32 dimensions, if it was compiled without alterations to the source code, we will only focus on lower-dimensional arrays for the purpose of illustration in this introduction.

Now, let us get started with NumPy by calling the array function to create a two-dimensional NumPy array, consisting of two rows and three columns, from a list of lists:

<sup>&</sup>lt;sup>6</sup>https://docs.scipy.org/doc/numpy/reference/generated/numpy.ndarray.html



<sup>&</sup>lt;sup>2</sup>http://www.labri.fr/perso/nrougier/from-python-to-numpy/

<sup>&</sup>lt;sup>3</sup>https://archive.org/details/NumPyBook

<sup>&</sup>lt;sup>4</sup>http://www.scipy-lectures.org/intro/numpy/index.html

<sup>&</sup>lt;sup>5</sup>https://docs.scipy.org/doc/numpy/reference/index.html

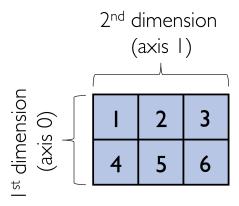


Figure F.1: Illustration of a 2-dimensional NumPy array.

By default, NumPy infers the type of the array upon construction. Since we passed Python integers to the array, the ndarray object aryld should be of type int 64 on a 64-bit machine, which we can confirm by accessing the dtype attribute:

```
1 >>> ary1d.dtype
2 dtype('int64')
```

If we want to construct NumPy arrays of different types, we can pass an argument to the dtype parameter of the array function; for example, dtype=np.int32 to create 32-bit arrays. For a full list of supported data types, please refer to the official NumPy documentation<sup>7</sup>. Once an array has been constructed, we can downcast or recast its type via the astype method as shown in the following example:

<sup>&</sup>lt;sup>7</sup>https://docs.scipy.org/doc/numpy/user/basics.types.html



In the following sections we will cover many more aspects of NumPy arrays; however, to conclude this basic introduction to the ndarray object, let us take a look at some of its handy attributes. For instance, the itemsize attribute returns the size of a single array element in bytes:

The code snippet above returned 8, which means that each element in the array (remember that ndarrays are homogeneous) takes up 8 bytes in memory. This result makes sense since the array ary2d has type int64 (64-bit integer), which we determined earlier, and 8 bits equals 1 byte. (Note that 'int64' is just a shorthand for np.int64.)

To return the number of elements in an array, we can use the size attribute, as shown below:

```
1 >>> ary2d.size
2 6
```

And the number of dimensions of an array (intuitively, you may think of *dimensions* as the *rank* of a tensor) can be obtained via the ndim attribute:

```
1 >>> ary2d.ndim
2 2
```

If we are interested in the number of elements along each array dimension (in the context of NumPy arrays, we may also refer to them as *axes*), we can access the shape attribute as shown below:

```
1 >>> ary2d.shape
2 (2, 3)
```



The shape is always a tuple; in the code example above, the two-dimensional ary object has two *rows* and *three* columns, (2, 3), if we think of it as a matrix representation. Similarly, the shape of the one-dimensional array only contains a single value:

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

Instead of passing lists or tuples to the array function, we can also provide a single float or integer, which will construct a zero-dimensional array (for instance, a representation of a scalar):

```
1  >>> scalar = np.array(5)
2  >>> scalar
3  array(5)
1  >>> scalar.ndim
2  0
1  >>> scalar.shape
2  ()
```

#### **F.2** Array Construction Routines

In the previous section, we used the array function to construct NumPy arrays from Python objects that are sequences or nested sequences – lists, tuples, nested lists, iterables, and so forth. While array is often our goto function for creating ndarray objects, NumPy implements a variety of functions for constructing arrays that may come in handy in different contexts. In this section, we will take a quick peek at those that we use most commonly – you can find a more comprehensive list in the official documentation<sup>8</sup>.

The array function works with most iterables in Python, including lists, tuples, and range objects; however, array does not support generator expressions. If we want parse generators directly, however, we can use the fromiter function as demonstrated below:

```
1 >>> def generator():
2 ... for i in range(10):
```

<sup>&</sup>lt;sup>8</sup>https://docs.scipy.org/doc/numpy/reference/routines.array-creation.html



```
if i % 2:

yield i

>>> gen = generator()

>>> np.fromiter(gen, dtype=int)

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

>>> # using 'comprehensions' the following

>>> # generator expression is equivalent to

>>> # the \code{generator} above

>>> generator_expression = (i for i in range(10) if i % 2)

>>> np.fromiter(generator_expression, dtype=int)

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

Next, we will take at two functions that let us create ndarrays of consisting of either ones and zeros by only specifying the elements along each axes (here: three rows and three columns):

Creating arrays of ones or zeros can also be useful as placeholder arrays, in cases where we do not want to use the initial values for computations but want to fill it with other values right away. If we do not need the initial values (for instance, '0.' or '1.'), there is also numpy.empty, which follows the same syntax as numpy.ones and np.zeros. However, instead of filling the array with a particular value, the empty function creates the array with non-sensical values from memory. We can think of zeros as a function that creates the array via empty and then sets all its values to 0. – in practice, a difference in speed is not noticeable, though.

NumPy also comes with functions to create identity matrices and diagonal matrices as ndarrays that can be useful in the context of linear algebra – a topic that we will explore later in this appendix.

```
1 >>> np.eye(3)
2 array([[ 1., 0., 0.],
```



Lastly, I want to mention two very useful functions for creating sequences of numbers within a specified range, namely, arange and linspace. NumPy's arange function follows the same syntax as Python's range objects: If two arguments are provided, the first argument represents the start value and the second value defines the stop value of a half-open interval:

```
1 >>> np.arange(4., 10.)
2 array([ 4., 5., 6., 7., 8., 9.])
```

Notice that arange also performs type inference similar to the array function. If we only provide a single function argument, the range object treats this number as the endpoint of the interval and starts at 0:

```
1 >>> np.arange(5)
2 array([0, 1, 2, 3, 4])
```

Similar to Python's range, a third argument can be provided to define the *step* (the default step size is 1). For example, we can obtain an array of all uneven values between one and ten as follows:

```
1 >>> np.arange(1., 11., 2)
2 array([ 1., 3., 5., 7., 9.])
```

The linspace function is especially useful if we want to create a particular number of evenly spaced values in a specified half-open interval:

```
1 >>> np.linspace(0., 1., num=5)
2 array([ 0. , 0.25, 0.5 , 0.75, 1. ])
```

#### F.3 Array Indexing

In this section, we will go over the basics of retrieving NumPy array elements via different indexing methods. Simple NumPy indexing and slicing



works similar to Python lists, which we will demonstrate in the following code snippet, where we retrieve the first element of a one-dimensional array:

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

Also, the same Python semantics apply to slicing operations. The following example shows how to fetch the first two elements in ary:

```
1 >>> ary[:2] # equivalent to ary[0:2]
2 array([1, 2])
```

If we work with arrays that have more than one dimension or axis, we separate our indexing or slicing operations by commas as shown in the series of examples below:

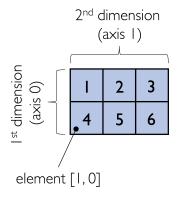


Figure F.2: NumPy array indexing.



#### F.4 Array Math and Universal Functions

In the previous sections, you learned how to create NumPy arrays and how to access different elements in an array. It is about time to cover one of the core features of NumPy that makes working with ndarray so efficient and convenient: vectorization. While we typically use for-loops if we want to perform arithmetic operations on sequence-like objects, NumPy provides vectorized wrappers for performing element-wise operations implicitly via so-called *ufuncs* – short for universal functions.

As of this writing, there are more than 60 ufuncs available in NumPy; ufuncs are implemented in compiled C code and very fast and efficient compared to vanilla Python. In this section, we will take a look at the most commonly used ufuncs, and I recommend you to check out the official documentation <sup>9</sup> for a complete list.

To provide an example of a simple usunc for element-wise addition, consider the following example, where we add a scalar (here: 1) to each element in a nested Python list:

```
1 >>> lst = [[1, 2, 3], [4, 5, 6]]
2 >>> for row_idx, row_val in enumerate(lst):
3 ... for col_idx, col_val in enumerate(row_val):
4 ... lst[row_idx][col_idx] += 1
5 >>> lst
6 [[2, 3, 4], [5, 6, 7]]
```

<sup>&</sup>lt;sup>9</sup>https://docs.scipy.org/doc/numpy/reference/ufuncs.html#available-ufuncs



This for-loop approach is very verbose, and we could achieve the same goal more elegantly using list comprehensions:

```
1 >>> lst = [[1, 2, 3], [4, 5, 6]]
2 >>> [[cell + 1 for cell in row] for row in lst]
3 [[2, 3, 4], [5, 6, 7]]
```

We can accomplish the same using NumPy's ufunc for element-wise scalar addition as shown below:

The ufuncs for basic arithmetic operations are add, subtract, divide, multiply, and exp (exponential). However, NumPy uses operator overloading so that we can use mathematical operators (+, -, /, \*, and \*\*) directly:

In the previous code snippets, we have seen examples of *binary* ufuncs, which are ufuncs that perform computations between two input arguments. In addition, NumPy implements several useful *unary* ufuncs that perform computations on a single array; examples are <a href="log">log</a> (natural logarithm), <a href="log">logarithm</a>), <a href="log">logarithm</a>),

Often, we want to compute the sum or product of array element along a given axis. For this purpose, we can use a ufunc's reduce operation. By default, reduce applies an operation along the first axis (axis=0). In the case of a two-dimensional array, we can think of the first axis as the rows of a matrix. Thus, adding up elements along rows yields the column sums of that matrix as shown in the following code snippet:



To compute the row sums of the array above, we can specify axis=1:

```
1 >>> np.add.reduce(ary, axis=1) # row sums
2 array([ 6, 15])
```

While it can be more intuitive to use reduce as a more general operation, NumPy also provides shorthands for specific operations such as product and sum. For example, sum (axis=0) is equivalent to add.reduce:

```
1 >>> ary.sum(axis=0) # column sums
2 array([5, 7, 9])
```

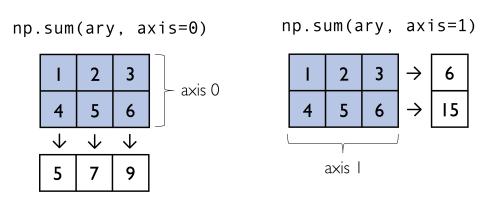


Figure F.3: The NumPy sum operation.

Note that np.sum(ary, ...) and ary.sum(...) are equivalent operations. As a word of caution, keep in mind that product and sum both compute the product or sum of the entire array if we do not specify an axis:

```
1 >>> ary.sum()
2 21
```

Other useful unary ufuncs are:

• mean (computes arithmetic average)



- std (computes the standard deviation)
- var (computes variance)
- np.sort (sorts an array)
- np.argsort (returns indices that would sort an array)
- np.min (returns the minimum value of an array)
- np.max (returns the maximum value of an array)
- np.argmin (returns the index of the minimum value)
- np.argmax (returns the index of the maximum value)
- array\_equal (checks if two arrays have the same shape and elements)

#### F.5 Broadcasting

A topic we glanced over in the previous section is broadcasting. Broadcasting allows us to perform vectorized operations between two arrays even if their dimensions do not match by creating implicit multidimensional grids. You already learned about ufuncs in the previous section where we performed element-wise addition between a scalar and a multidimensional array, which is just one example of broadcasting.

Figure F.4: An illustration of broadcasting showing the addition of a single number to every element in an array.

Naturally, we can also perform element-wise operations between arrays of equal dimensions:



```
1 >>> ary1 = np.array([1, 2, 3])
2 >>> ary2 = np.array([4, 5, 6])
3 >>> ary1 + ary2
4 array([5, 7, 9])
```

In contrast to what we are used from linear algebra, we can also add arrays of different shapes. In the example above, we will add a one-dimensional to a two-dimensional array, where NumPy creates an implicit multidimensional grid from the one-dimensional array ary1:

Figure F.5: An illustration of broadcasting showing the addition of a 1D to a 2D NumPy array.

Keep in mind though that the number of elements along the explicit axes and the implicit grid have to match so that NumPy can perform a sensical operation. For instance, the following example should raise a ValueError, because NumPy attempts to add the elements from the first axis of the left array (2 elements) to the first axis of the right array (3 elements):

```
1 >>> try:
2 ... ary3 + np.array([1, 2])
3 >>> except ValueError as e:
4 ... print('ValueError:', e)
5 ValueError: operands could not be broadcast
6 together with shapes (2,3) (2,)
```



So, if we want to add the 2-element array to the columns in ary3, the 2-element array must have two elements along its *first* axis as well:

## F.6 Advanced Indexing – Memory Views and Copies

In the previous sections, we have used basic indexing and slicing routines. It is important to note that basic integer-based indexing and slicing create so-called *views* of NumPy arrays in memory. In other words, via a view, we are accessing the same array but just "view" it differently.

Working with views can be highly desirable since it avoids making unnecessary copies of arrays to save memory resources. To illustrate the concept of memory views, let us walk through a simple example where we access the first row in an array, assign it to a variable, and modify that variable:

As we can see in the example above, changing the value of first\_row also affected the original array. The reason for this is that ary[0] created a view of the first row in ary, and its elements were then incremented by 99. The same concept applies to slicing operations:



If we are working with NumPy arrays, it is always important to be aware that slicing creates views – sometimes it is desirable since it can speed up our code and save computational resources by avoiding to create unnecessary copies in memory. However, in certain scenarios we want force a copy of an array; we can do this via the copy method as shown below:

One way to check if two arrays might share memory is to use the NumPy's may\_share\_memory function. However, be aware that it is a heuristic that can return false negatives or false positives in rare cases. The next code snippet shows an example of may\_share\_memory applied to a view (first\_row) and copy (second\_row) of the array elements from ary:



In addition to basic single-integer indexing and slicing operations, NumPy supports advanced indexing routines called *fancy* indexing. Via fancy indexing, we can use tuple or list objects of non-contiguous integer indices to return desired array elements. Since fancy indexing can be performed with non-contiguous sequences, it cannot return a view – a contiguous slice from memory. Thus, fancy indexing always returns a copy of an array – it is important to keep that in mind. The following code snippets show some fancy indexing examples:

Finally, we can also use Boolean masks for indexing – that is, arrays of True and False values. Consider the following example, where we return all values in the array that are greater than 3:



Using a *mask*, we can select elements given our desired criteria:

```
1 >>> ary[greater3_mask]
2 array([4, 5, 6])
```

We can also chain different selection criteria using the logical *and* operator "&" or the logical *or* operator "|." The example below demonstrates how we can select array elements that are greater than 3 *and* divisible by 2:

```
1 >>> ary[(ary > 3) & (ary % 2 == 0)]
2 array([4, 6])
```

Note that indexing using Boolean arrays (i.e., masks) is also considered "fancy indexing" and thus always returns a copy of the array and not view.

#### F.7 Comparison Operators and Mask

In the previous section, we already briefly introduced the concept of Boolean masks in NumPy. Boolean masks are bool-type arrays (storing True and False values) that have the same shape as a certain target array. For example, consider the following 4-element array in the next code snippet. Using comparison operators (such as <, >, <=, and >=), we can create a Boolean mask of that array which consists of True and False elements depending on whether a condition is met in the target array (here: ary):

```
1 >>> ary = np.array([1, 2, 3, 4])
2 >>> mask = ary > 2
3 >>> mask
4 array([False, False, True, True], dtype=bool)
```

One we created such a Boolean mask, we can use it to select certain entries from the target array – those entries that match the condition upon which the mask was created:

```
1 >>> ary[mask]
2 >>> array([3, 4])
```



Beyond the selection of elements from an array, Boolean masks can also come in handy when we want to count how many elements in an array meet a certain condition:

```
1 >>> mask
2 array([False, False, True, True], dtype=bool)
3 >>> mask.sum()
4 2
```

If we are interested in the index positions of array elements that meet a certain condition, we can use the nonzero() method on a mask, as follows:

```
1 >>> mask.nonzero()
2 (array([2, 3]),)
```

Note that selecting index elements in the way suggested above is a two-step process of creating a mask and then applying the non-zero method:

```
1 >>> (ary > 2).nonzero()
2 (array([2, 3]),)
```

An alternative approach to the index selection by a condition is using the np.where method:

```
1 >>> np.where(ary > 2)
2 (array([2, 3]),)
1 >>> np.where(ary > 2, 1, 0)
2 array([0, 0, 1, 1])
```

Notice that we use the np. where function with three arguments:

```
np.where (condition, x, y), which is interpreted as
```

If condition is True, yield x, otherwise yield y.

Or more concretely, what we have done in the previous example is to assign 1 to all elements greater than 2, and 0 to all other elements. Of course, this can also be achieved by using Boolean masks "manually:"

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



```
3 >>> ary[mask] = 1
4 >>> ary[~mask] = 0
5 >>> ary
6 array([0, 0, 1, 1])
```

The ~ operator in the example above is one of the logical operators in NumPy:

```
• A: & or np.bitwise_and
```

- Or: | or np.bitwise\_or
- Xor: ^ or np.bitwise\_xor
- Not: ~ or np.bitwise\_not

Logical operators allow us to chain an arbitrary number of conditions to create even more "complex" Boolean masks. For example, using the "Or" operator, we can select all elements that are greater than 3 or smaller than 2 as follows:

```
1 >>> ary = np.array([1, 2, 3, 4])
2 >>> (ary > 3) | (ary < 2)
3 array([ True, False, False, True], dtype=bool)</pre>
```

And, for example, to negate the condition, we can use the ~ operator:

```
1 >>> ~((ary > 3) | (ary < 2))
2 array([False, True, True, False], dtype=bool)</pre>
```

#### F.8 Random Number Generators

In machine learning and deep learning, we often have to generate arrays of random numbers – for example, the initial values of our model parameters before optimization. NumPy has a random subpackage to create random numbers and samples from a variety of distributions conveniently. Again, I encourage you to browse through the more comprehensive numpy.random documentation<sup>10</sup> for a more comprehensive list of functions for random sampling. To provide a brief overview of the pseudo-random number generators that we will use most commonly, let's start with drawing a random sample from a uniform distribution:

<sup>&</sup>lt;sup>10</sup>https://docs.scipy.org/doc/numpy/reference/routines.random.html



```
1 >>> np.random.seed(123)
2 >>> np.random.rand(3)
3 array([ 0.69646919,  0.28613933,  0.22685145])
```

In the code snippet above, we first seeded NumPy's random number generator. Then, we drew three random samples from a uniform distribution via random.rand in the half-open interval [0, 1). I highly recommend the seeding step in practical applications as well as in research projects, since it ensures that our results are reproducible. If we run our code sequentially – for example, if we execute a Python script – it should be sufficient to seed the random number generator only once at the beginning to enforce reproducible outcomes between different runs. However, it is often useful to create separate RandomState objects for various parts of our code, so that we can test methods of functions reliably in unit tests.

Working with multiple, separate <code>RandomState</code> objects can also be useful if we run our code in non-sequential order; for example, if we are experimenting with our code in interactive sessions or Jupyter Notebook environments. The example below shows how we can use a <code>RandomState</code> object to create the same results that we obtained via <code>np.random.rand</code> in the previous code snippet:

```
1 >>> rng1 = np.random.RandomState(seed=123)
2 >>> rng1.rand(3)
3 array([ 0.69646919,  0.28613933,  0.22685145])
```

Another useful function that we will often use in practice is randn, which returns a random sample of floats from a standard normal distribution  $N(\mu, \sigma^2)$ , where the mean,  $(\mu)$  is zero and unit variance  $(\sigma = 1)$ . The example below creates a two-dimensional array of such z-scores:

```
1 >>> rng2 = np.random.RandomState(seed=123)
2 >>> z_scores = rng2.randn(100, 2)
```

NumPy's random functions rand and randn take an arbitrary number of integer arguments, where each integer argument specifies the number of elements along a given axis – the z\_scores array should now refer to an array of 100 rows and two columns. Let us now visualize (Figure F.6) how our random sample looks like using Matplotlib<sup>11</sup> [Barrett et al., 2005, Hunter, 2007]:

<sup>&</sup>lt;sup>11</sup>https://matplotlib.org



```
1 >>> import matplotlib.pyplot as plt
2 >>> plt.scatter(z_scores[:, 0], z_scores[:, 1])
3 >>> plt.show()
```

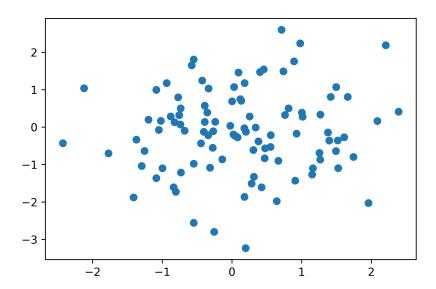


Figure F.6: 100 random data points samples from a standard-normal distribution via np.random.randn.

If we want to draw a random sample from a non-standard normal distribution, we can simply add a scalar value to the array elements to shift the mean of the sample, and we can multiply the sample by a scalar to change its standard deviation. The following code snippet will change the properties of our random sample as if it has been drawn from a normal distribution N(5,4):

```
1 >>> rng3 = np.random.RandomState(seed=123)
2 >>> scores = 2. * rng3.randn(100, 2) + 5.
3 >>> plt.scatter(scores[:, 0], scores[:, 1])
4 >>> plt.show()
```



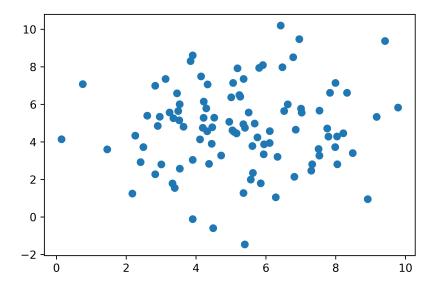


Figure F.7: 100 random data points generated via np.random.rand.

Note that in the example above, we multiplied the z-scores by a standard deviation of 2 – the standard deviation of a sample is the square root of the variance  $\sigma^2$ . Also, notice that all elements in the array were updated when we multiplied it by a scalar or added a scalar. In the next section, we will discuss NumPy's capabilities regarding such element-wise operations in more detail.

## F.9 Reshaping Arrays

In practice, we often run into situations where existing arrays do not have the *right* shape to perform certain computations. As you might remember from the beginning of this appendix, the size of NumPy arrays is fixed. Fortunately, this does not mean that we have to create new arrays and copy values from the old array to the new one if we want arrays of different shapes – the size is fixed, but the shape is not. NumPy provides a reshape method that allow us to obtain a view of an array with a different shape. For example, we can reshape a one-dimensional array into a two-dimensional one using reshape as follows:



While we need to specify the desired elements along each axis, we need to make sure that the reshaped array has the same number of elements as the original one. However, we do not need to specify the number elements in each axis; NumPy is smart enough to figure out how many elements to put along an axis if only one axis is unspecified (by using the placeholder -1):

We can, of course, also use reshape to flatten an array:

Note that NumPy also has a shorthand for that called ravel:

```
1 >>> ary2d.ravel()
2 array([1, 2, 3, 4, 5, 6])
```

A function related to ravel is flatten. In contrast to ravel, flatten returns a copy, though:

```
1 >>> np.may_share_memory(ary2d.flatten(), ary2d)
2 False
```



```
1 >>> np.may_share_memory(ary2d.ravel(), ary2d)
2 True
```

Sometimes, we are interested in merging different arrays. Unfortunately, there is no efficient way to do this without creating a new array, since NumPy arrays have a fixed size. While combining arrays should be avoided if possible – for reasons of computational efficiency – it is sometimes necessary. To combine two or more array objects, we can use NumPy's concatenate function as shown in the following examples:

#### F.10 Linear Algebra with NumPy Arrays

Most of the operations in machine learning and deep learning are based on concepts from linear algebra. In this section, we will take a look how to perform basic linear algebra operations using NumPy arrays.

#### Infobox F.10.1 Numpy Matrix Objects

I want to mention that there is also a special  $\mathtt{matrix}^a$  type in NumPy. NumPy  $\mathtt{matrix}$  objects are analogous to NumPy arrays but are restricted to two dimensions. Also, matrices define certain operations differently than arrays; for instance, the  $\star$  operator performs matrix multiplication instead of element-wise multiplication. However, NumPy  $\mathtt{matrix}$  is less popular in the science community compared to the more general array data structure. Since we are also going to work



with arrays that have more than two dimensions (for example, when we are working with convolutional neural networks), we will not use NumPy matrix data structures in this book.

<sup>a</sup>https://docs.scipy.org/doc/numpy/reference/generated/numpy.matrix.html

Intuitively, we can think of one-dimensional NumPy arrays as data structures that represent row vectors:

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

Similarly, we can use two-dimensional arrays to create column vectors:

Instead of reshaping a one-dimensional array into a two-dimensional one, we can simply add a new axis as shown below:

Note that in this context, np.newaxis behaves like None:

All three approaches listed above, using reshape (-1, 1), np.newaxis, or None yield the same results – all three approaches create views not copies of the row\_vector array.

As we remember from the Linear Algebra appendix, we can think of a vector as a  $1 \times n$ -dimensional matrix (row vector) or  $n \times 1$ -dimensional matrix (column vector). Then, to perform matrix multiplication between



matrices, we learned that number of columns of the left matrix must match the number of rows of the matrix to the right. In NumPy, we can perform matrix multiplication via the matmul function:

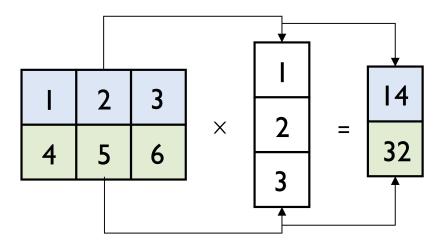


Figure F.8: PLACEHOLDER

However, if we are working with matrices and vectors, NumPy can be quite forgiving if the dimensions of matrices and one-dimensional arrays do not match exactly – thanks to broadcasting. The following example yields the same result as the matrix-column vector multiplication, except that it returns a one-dimensional array instead of a two-dimensional one:

```
1 >>> np.matmul(matrix, row_vector)
2 array([14, 32])
```

Similarly, we can compute the dot-product between two vectors (here: the vector norm)



```
1 >>> np.matmul(row_vector, row_vector)
2 14
```

NumPy has a special dot function that behaves similar to matmul on pairs of one- or two-dimensional arrays – its underlying implementation is different though, and one or the other can be slightly faster on specific machines and versions of BLAS <sup>12</sup>:

```
1 >>> np.dot(row_vector, row_vector)
2 14

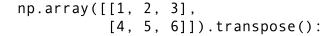
1 >>> np.dot(matrix, row_vector)
2 array([14, 32])

1 >>> np.dot(matrix, column_vector)
2 array([[14],
3 [32]])
```

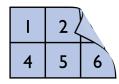
Similar to the examples above we can use matmul or dot to multiply two matrices (here: two-dimensional arrays). In this context, NumPy arrays have a handy transpose method (Figure F.9) to transpose matrices if necessary:



<sup>&</sup>lt;sup>12</sup>https://en.wikipedia.org/wiki/Basic\_Linear\_Algebra\_Subprograms



ı	2	3
4	5	6





I	4
2	5
3	6

Figure F.9: Illustration of NumPy's transpose method.

While transpose can be annoyingly verbose for implementing linear algebra operations – think of PEP8's  $^{13}$  80 character per line recommendation – NumPy has a shorthand for that, T:



<sup>&</sup>lt;sup>13</sup>https://www.python.org/dev/peps/pep-0008/

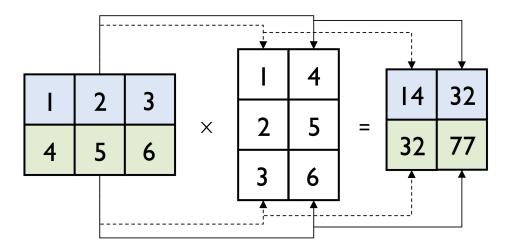


Figure F.10: Multiplication of a matrix with its transpose.

While this section demonstrates some of the basic linear algebra operations carried out on NumPy arrays that we use in practice, you can find an additional function in the documentation of NumPy's submodule for linear algebra: numpy.linalg<sup>14</sup>. If you want to perform a particular linear algebra routine that is not implemented in NumPy, it is also worth consulting the scipy.linalg documentation <sup>15</sup> – SciPy [Jones et al., 2018] is a library for scientific computing built on top of NumPy.

## F.11 Set Operations

Appendix B: Algebra Basics introduced the basics behind set theory, which are visually summarized in Figure F.11, and NumPy implements several functions that allow us to work with sets efficiently so that we are not restricted to Python's built-in functions (and converting back and forth between Python sets and NumPy arrays).



<sup>&</sup>lt;sup>14</sup>https://docs.scipy.org/doc/numpy/reference/routines.linalg.html

<sup>&</sup>lt;sup>15</sup>https://docs.scipy.org/doc/scipy/reference/linalg.html

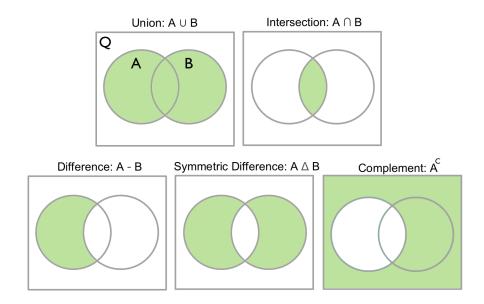


Figure F.11: Common set operations.

Remember that a set is essentially a collection of unique elements. Given an array, we can generate such a "set" using the np.unique function:

```
1 >>> ary = np.array([1, 1, 2, 3, 1, 5])
2 >>> ary_set = np.unique(ary)
3 >>> ary_set
4 >>> array([1, 2, 3, 5])
```

However, we have to keep in mind that unlike Python sets, the output of np.unique is a regular NumPy array, not specialized data structure that does not allow for duplicate entries. The set operations for example, set union (np.union1d), set difference (np.setdiff1d), or set intersection (np.intersect1d) would return the same results whether array elements are unique or not. However, setting their optional assume\_unique argument can speed up the computation.

```
1 >>> ary1 = np.array([1, 2, 3])
2 >>> ary2 = np.array([3, 4, 5, 6])
3 >>> np.intersect1d(ary1, ary2, assume_unique=True)
4 array([3])
```



```
1 >>> np.setdiffld(ary1, ary2, assume_unique=True)
2 array([1, 2])
1 >>> np.union1d(ary1, ary2) # does not have assume_unique
2 array([1, 2, 3, 4, 5, 6])
```

Note that NumPy does not have a function for the symmetric set difference, but it can be readily computed by composition:

```
1 >>> np.union1d(np.setdiff1d(ary1, ary2,
2 ... assume_unique=True),
3 ... np.setdiff1d(ary2, ary1,
4 ... assume_unique=True))
5 array([1, 2, 4, 5, 6])
```

For a complete list of set operations in NumPy, please see the official Set Routines<sup>16</sup> documentation.

#### F.12 Serializing NumPy Arrays

In the context of computer science, serialization is a term that refers to storing data or objects in a different format that can be used for reconstruction later. For example, in Python, we can use the pickle library to write Python objects as bytecode (or binary code) to a local drive. NumPy offers a data storage format (NPy<sup>17</sup>) that is especially well-suited (compared to regular pickle files) for storing array data. There are three different functions that can be used for this purpose:

```
np.savenp.saveznp.savez_compressed
```

Starting with np.save, this function saves a single array to a so-called .npy file:

```
1 >>> ary1 = np.array([1, 2, 3])
2 >>> np.save('ary-data.npy', ary1)
3 >>> np.load('ary-data.npy')
4 array([1, 2, 3])
```

 <sup>16</sup>https://docs.scipy.org/doc/numpy/reference/routines.set.html
 17https://docs.scipy.org/doc/numpy/neps/npy-format.html



Suffice it to say, the the np.load a universal function for loading data that has been serialized via NumPy back into the current Python session.

The np.savez is slightly more powerful than the np.save function as it generates an archive consisting of 1 or more .npy files and thus allows us to save multiple arrays at once.

```
np.savez(<mark>'ary-data.npz',</mark> ary1, ary2)
```

However, note that if the np.load function will wrap the arrays as a NpzFile object from which the individual arrays can be accessed via keys like values in a Python dictionary:

```
1 >>> d = np.load('ary-data.npz')
2 >>> d.keys()
3 ['arr_0', 'arr_1']
1 >>> d['arr_0']
2 array([1, 2, 3])
```

As demonstrated above, the arrays are numerated in ascending order with an arr\_ prefix ('arr\_0', 'arr\_1', etc.). To use a custom naming convention, we can simply provide keyword arguments:

#### F.13 Conclusion

We have covered a lot of material in this appendix. If you are new to NumPy, its functionality can be quite overwhelming at first. Good news is that we do not need to master all the different concepts at once before we can get started using NumPy in our applications. In my opinion, the most useful yet most difficult concepts are NumPy's broadcasting rules and to distinguish between views and copies of arrays. However, with some experimentation, you can quickly get the hang of it and be able to write elegant and efficient NumPy code.



## **Bibliography**

- [Barrett et al., 2005] Barrett, P., Hunter, J., Miller, J. T., Hsu, J.-C., and Greenfield, P. (2005). matplotlib–a portable python plotting package. In *Astronomical Data Analysis Software and Systems XIV*, volume 347, page 91.
- [Hunter, 2007] Hunter, J. D. (2007). Matplotlib: A 2d graphics environment. *Computing in science & engineering*, 9(3):90–95.
- [Jones et al., 2018] Jones, E., Oliphant, T., Peterson, P., et al. (2001–2018). SciPy: Open source scientific tools for Python. [Online; accessed 05/30/2018].
- [Oliphant, 2006] Oliphant, T. E. (2006). *A guide to NumPy*, volume 1. Trelgol Publishing USA.
- [Walt et al., 2011] Walt, S. v. d., Colbert, S. C., and Varoquaux, G. (2011). The numpy array: a structure for efficient numerical computation. *Computing in Science & Engineering*, 13(2):22–30.

# **Abbreviations and Terms**

**CNN** [Convolutional Neural Network]



## Index

