

# Getting Started with NumPy

- Introduction to Python data science
- NumPy arrays
- Manipulating array elements
- Manipulating array shape

# Introduction to Python data science

- Overview
- Python libraries for data science
- Getting the data science libraries

# Overview

Python is a popular choice for data science and machine learning

Attractive characteristics of Python:

- Dynamic language, so it's good for rapid exploratory coding
- Relatively simple syntax, so it's easier to become proficient
- Popular in schools and universities, so the skills are out there!

# Python libraries for data science

NumPy is a numeric processing API for Python

- Fast mathematical computation of numeric arrays and matrices

Pandas provides additional features based on NumPy

- Additional support for indexing, reading/writing CSV/Excel, etc.

Matplotlib is a graphical plotting API for Python

- Similar to Matlab, allows you to plot graphs, charts, etc.

Scikit-Learn is a machine learning library for Python

- Implements many supervised/unsupervised learning algorithms

# Getting the data science libraries

If you're using Anaconda, these are already downloaded for you.

If you're using a standalone Python distribution, you'll need to install the libraries you need.

```
1  pip install numpy
2
3  pip install openpyxl
4  pip install xlrd
5  pip install matplotlib
6
7  pip install pandas
```

# NumPy arrays

- Getting Started with NumPy arrays
- Techniques for creating NumPy arrays
- Reading CSV data
- Visualizing data

# Getting Started with NumPy arrays (1/2)

NumPy holds data in N-dimensional arrays

- An array is an instance of the `numpy.ndarray` class
- <https://docs.scipy.org/doc/numpy/reference/generated/numpy.ndarray.html>

All the data in a NumPy array is the same type

- This allows NumPy to store and process the data efficiently
- This is a very important point!

Why are NumPy arrays more efficient than Python lists?

- Python is dynamically typed, so every object contains metadata that identifies the type at run time
- In a Python list, every item contains this metadata - eek!
- In a NumPy array, only the array itself contains the metadata



# Getting Started with NumPy arrays (2/2)

## Example

- Creates a NumPy array from a Python list
- Gets the shape of the array, via the shape property
- Gets the data type of array elements, via the dtype property

```
1  # Import the NumPy module.
2  import numpy as np
3
4  # Create a 1D NumPy array from a Python list.
5  a = np.array([1, 2, 3])
6  print('Data values in a\n', a)
7  print('Shape of a:', a.shape)
8  print('Data type in a:', a.dtype)
9
10 # Create a 2D Numpy array from a Python list of lists.
11 b = np.array([[1, 2, 3], [4, 5, 6]])
12 print('\nData values in b\n', b)
13 print('Shape of b:', b.shape)
14 print('Data type in b:', b.dtype)
```

For a full list of NumPy standard types, see:

# Techniques for creating NumPy arrays (1/2)

There are lots of ways to create a NumPy array

```
1  import numpy as np
2
3  # Create array with mixed types - NumPy converts elements "upwards".
4  a = np.array([1, 2, 3.14])
5
6  # Create array with a specified type.
7  b = np.array([1, 2, 3], dtype='float64')
8
9  # Create array from a numeric range.
10 c = np.arange(0, 20, 2)
11
12 # Create array of elements, linear spaced.
13 d = np.linspace(0.0, 1.0, 11)
14
15 # Create array of zeros.
16 # You can specify multi-dim arrays too.
17 e = np.zeros(5)
18
19 # Create array of ones.
20 f = np.ones(5)
21
22 # Create array of elements, with specified value.
23 g = np.full(5, 1.22)
```

# Techniques for creating NumPy arrays (2/2)

You can also create random arrays, which can be handy

```
1  # Import the NumPy module.
2  import numpy as np
3
4  # Create random values in range [0.0, 1.0).
5  a = np.random.random(10)
6
7  # Create normally-distributed random values.
8  b = np.random.normal(5, 2, 10)
9
10 # Create random integers in range [0, 101).
11 c = np.random.randint(0, 101, 10)
```

# Reading CSV data

A common requirement is to read data from a CSV file

- The easiest way to do this is via the Pandas `read_csv()` function

Pandas reads values into a multi-column DataFrame

- You can then extract a column into a NumPy array

```
1  import numpy as np
2  import pandas as pd
3
4  # Read a csv file, get a Pandas DataFrame back.
5  dataframe = pd.read_csv('WorldCupWinners.csv')
6
7  # Get the 'Teams' column.
8  teams = np.array(dataframe['Team'])
9  print(teams)
```

We'll not have time to dive into Pandas on this course but it's a topic worth exploring!

# Visualizing data (1/2)

Visualization is an important aid to help you understand the shape and meaning of data

You can use the Matplotlib library to visualize data in lots of different ways

- Line graphs
- Scatter graphs
- Bar-charts
- Pie-charts
- Histograms
- Etc.

# Visualizing data (2/2)

Here's a simple example of how to visualize data using Matplotlib - we'll see more plotting features later

```
1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4
5  data = np.array([1, 19, 76, 45, 34, 42, 30, 5, 77, 54, 89])
6
7  plt.xlabel('Element in array')
8  plt.ylabel('Value')
9  plt.plot(data)
10 plt.show()
```

# Manipulating array elements

- Indexing into an array
- Slicing an array
- Accessing a specific column or row
- Aside: Views vs. copies

# Indexing into an array

Indexing into a NumPy array is quite intuitive

- `[i]` Access element from start, first element is at `[0]`
- `[-i]` Access element from end, last element is at `[-1]`
- `[r,c]` Access element in 2-D array (etc. for higher dimensions)

```
1  import numpy as np
2
3  # Create a 1-D array, index into it, and modify elements.
4  a = np.array([0, 10, 20, 30, 40, 50, 60, 70])
5  print(a)
6  print(a[1])          # 10
7  print(a[-1])         # 70
8  a[1] = 111
9  print(a)
10
11 # Create a 2-D array, index into it, and modify elements.
12 b = np.array([[0, 10, 20, 40], [50, 60, 70, 80]])
13 print(b)
14 print(b[0, 1])       # 10
15 print(b[0, -1])      # 40
16 print(b[-1, 1])      # 60
17 print(b[-1, -1])     # 80
18 b[0, 1] = 111
```



# Slicing an array

You can slice into an array using a [start:stop:step] index

- start Default start is 0
- stop Default stop is the size of the dimension
- step Default step is 1

```
1  import numpy as np
2
3  # Create a 1-D array, and get various slices.
4  a = np.array([0, 10, 20, 30, 40, 50, 60, 70, 80])
5  print(a[3:])          # [30 40 50 60 70 80]
6  print(a[:3])          # [ 0 10 20]
7  print(a[3:7])         # [30 40 50 60]
8  print(a[3:7:2])       # [30 50]
9  print(a[3::2])        # [30 50 70]
10 print(a[::2])         # [ 0 20 40 60 80]
11
12 # Create a 2-D array, and get various slices in each dimension.
13 b = np.array([[0, 10, 20], [30, 40, 50], [60, 70, 80]])
14 print(b[1:, 1:])       # [[40 50] [70 80] ]
15 print(b[:2, :2])      # [[ 0 10] [30 40] ]
16 print(b[::2, ::2])     # [[ 0 20] [60 80] ]
17
```

# Accessing a specific column or row

To get a specific column or row in a multidimension array:

- Use an empty slice to skip a dimension
- E.g. in a 2D array, `[:,1]` gets column 1
- E.g. in a 2D array, `[1,:]` gets row 1

```
1  import numpy as np
2
3  a = np.array([[0, 10, 20], [30, 40, 50], [60, 70, 80]])
4
5  # To access a specific column ...
6  print(a[:, 0])    # [ 0 30 60]
7  print(a[:, 1])    # [10 40 70]
8  print(a[:, 2])    # [20 50 80]
9
10 # To access a specific row ...
11 print(a[0, :])    # [ 0 10 20]
12 print(a[1, :])    # [30 40 50]
13 print(a[2, :])    # [60 70 80]
14
15 # To access a specific row, simpler syntax ...
16 print(a[0])       # [ 0 10 20]
17 print(a[1])       # [30 40 50]
```

# Aside: Views vs. copies

When you get an array slice/row/column, you get a view on the data

- If you make any changes, it will change the actual data

If you want to get a copy of the data:

- Call `copy()` on the slice/row/column

```
1  import numpy as np
2
3  # Demonstrate views.
4  a = np.array([[0, 10, 20], [30, 40, 50], [60, 70, 80]])
5  col0View = a[:, 0]
6  col0View[2] = 600
7  print(col0View)
8  print(a)      # [[ 0  10  20] [ 30  40  50] [600  70  80]]
9
10 # Demonstrate copies.
11 b = np.array([[0, 10, 20], [30, 40, 50], [60, 70, 80]])
12 col0Copy = a[:, 0].copy()
13 col0Copy[2] = 600
14 print(col0Copy)
15 print(b)      # [[ 0  10  20] [ 30  40  50] [60  70  80]]
```

# Manipulating array shape

- Reshaping an array
- Creating new axes
- Concatenating arrays
- Stacking arrays vertically or horizontally
- Splitting an array

# Reshaping an array

Reshaping is a simple and common technique for creating multidimensional arrays

- Create a 1D array initially (typically)
- Reshape it to a multidimensional array (must be compatible shape)
- The multidimensional array is a view onto the original 1D array

```
1  import numpy as np
2
3  # Create 1D array initially, for simplicity.
4  a = np.arange(9)
5  print(a)           # [0 1 2 3 4 5 6 7 8]
6  print(a.shape)     # (9,)
7
8  # Reshape as 2D array (view on a).
9  b = a.reshape((3,3))
10 print(b)           # [[0 1 2] [3 4 5] [6 7 8]]
11 print(b.shape)     # (3, 3)
12
13 # Changing items in b will change values in underlying a.
14 b[0,0] = 99
15 print(a)           # [99  1  2  3  4  5  6  7  8]
16 print(b)           # [[99  1  2] [ 3  4  5] [ 6  7  8]]
```

# Creating new axes

Another useful technique is create new axes for an array

- Create a 1D array initially (typically)
- Create a new column or row, using `np.newaxis`

```
1  import numpy as np
2
3  # Create 1D array initially, for simplicity.
4  a = np.arange(5)
5  print(a)           # [0 1 2 3 4]
6  print(a.shape)     # (5,)
7
8  # Create 2D array with 1 row, 5 columns.
9  b = a[np.newaxis, :]
10 print(b)           # [[0 1 2 3 4]]
11 print(b.shape)     # (1, 5)
12
13 # Create 2D array with 5 rows, 1 column.
14 c = a[:, np.newaxis]
15 print(c)           # [[0] [1] [2] [3] [4]]
16 print(c.shape)     # (5, 1)
```

# Concatenating arrays (1/2)

You can concatenate same-size arrays together

- `np.concatenate()` - you can specify the axis to concatenate on

Here's a simple example that concatenates 1D arrays

```
1  import numpy as np
2
3  # Create some 1D arrays.
4  a = np.array([ 0,  1])
5  b = np.array([10, 11])
6  c = np.array([20, 21])
7
8  # Concatenate the 1D arrays.
9  result = np.concatenate([a, b, c])
10 print(result)          # [0 1 10 11 20 21]
11 print(result.shape)    # (6,)
```

# Concatenating arrays (2/2)

Here's an example that concatenates 2D arrays

- Note the optional axis parameter (default is 0)

```
1  import numpy as np
2
3  # Create some 2D arrays.
4  a = np.array([[ 0,  1], [10, 11]])
5  b = np.array([[20, 21], [30, 31]])
6  c = np.array([[40, 41], [50, 51]])
7
8  # Concatenate on axis 0 (this is the default, so can omit axis parameter).
9  result1 = np.concatenate([a, b, c], axis=0)
10 print(result1)          # [[0 1] [10 11] [20 21] [30 31] [40 41] [50 51]]
11 print(result1.shape)    # (6, 2)
12
13 # Concatenate on axis 1.
14 result2 = np.concatenate([a, b, c], axis=1)
15 print(result2)          # [[0 1 20 21 40 41] [10 11 30 31 50 51]]
16 print(result2.shape)    # (2, 6)
```



# Stacking arrays vertically or horizontally

You can stack different-size arrays together

- `np.vstack()` - stack vertically (must have same no. of cols)
- `np.hstack()` - stack horizontally (must have same no. of rows)

```
1  import numpy as np
2
3  # Create some arrays with same number of columns (2), and stack vertically.
4  a = np.array([10, 11])
5  b = np.array([[20, 21], [30, 31]])
6  result1 = np.vstack([a, b])
7  print(result1)          # [[10 11] [20 21] [30 31]]
8  print(result1.shape)    # (3, 2)
9
10 # Create some arrays with same number of rows (2), and stack horizontally.
11 c = np.array([[40, 41], [50, 51]])
12 d = np.array([[60], [61]])
13 result2 = np.hstack([c, d])
14 print(result2)          # [[40 41 60] [50 51 61]]
15 print(result2.shape)    # (2, 3)
```

# Splitting an array

You can split an array into subarrays

- `np.split()`
- `np.vsplit()`
- `np.hsplit()`

```
1  import numpy as np
2
3  # Split a 1D array.
4  a = np.arange(16)
5  a1, a2, a3, a4 = np.split(a, [2, 5, 9])
6  print('\na1\n', a1)    # [0 1]
7  print('\na2\n', a2)    # [2 3 4]
8  print('\na3\n', a3)    # [5 6 7 8]
9  print('\na4\n', a4)    # [9 10 11 12 13 14 15]
10
11 # Split a 2D vertically.
12 b = np.arange(16).reshape((4, 4))
13 b1, b2 = np.vsplit(b, [3])
14 print('\ntop\n', b1)    # [[0 1 2 3] [4 5 6 7] [8 9 10 11]]
15 print('\nbottom\n', b2) # [[12 13 14 15]]
16
17 # Split a 2D horizontally.
18 a = np.arange(16).reshape((4, 4))
```

Any questions?