



Lab 2. Fundamental Python libraries for Data Scientists: NumPy and SciPy, Scikit-Learn, Matplotlib and Seaborn

Learning Objectives

Import and run code examples using fundamental Python libraries for Data Scientists :

- NumPy and SciPy
- Scikit-Learn
- Matplotlib and Seaborn

Evaluation - 3% of total course mark

Instructions

1. Read Lab material
2. Run code examples from Lab Reading Material either in Jupyter or VS Code
3. Pick five interesting Numpy array functions and three interesting SciPy functions by going through the documentation. Run and **modify** the starter code to illustrate their usage with interesting example.
4. Add LinearDiscriminantAnalysis, KNeighborsClassifier, GaussianNB algorithms into scikit-learn example below. Comment on which algorithm demonstrates the best performance.
5. Use Jupyter notebook for this lab and include code, text, your comments, observations, and visualizations. Please make sure that the submitted notebooks have been run and the cell outputs are visible. Once created, submit it in pdf format.

Scoring Rubric

Category	Criteria	Maximum (Points)
Lab Examples	All python code is exceptionally well organized and very easy to follow.	0.5
Numpy, SciPy Examples	Additional five Numpy functions and three SciPy functions used properly.	1
scikit-learn examples	LinearDiscriminantAnalysis, KNeighborsClassifier, GaussianNB algorithms properly used to build a model. Comments provided on the performance of the best algorithm.	1.5



Lab 2 Material

Basic data types

Like most languages, Python has a number of basic types including integers, floats, booleans, and strings. These data types behave in ways that are familiar from other programming languages.

Numbers

Integers and floats work as you would expect from other languages. Note that unlike many languages, Python does not have unary increment ($x++$) or decrement ($x--$) operators.

```
x = 3
print(type(x)) # Prints "<class 'int'>"
print(x)       # Prints "3"
print(x + 1)   # Addition; prints "4"
print(x - 1)   # Subtraction; prints "2"
print(x * 2)   # Multiplication; prints "6"
print(x ** 2)  # Exponentiation; prints "9"
x += 1
print(x)       # Prints "4"
x *= 2
print(x)       # Prints "8"
y = 2.5
print(type(y)) # Prints "<class 'float'>"
print(y, y + 1, y * 2, y ** 2) # Prints "2.5 3.5 5.0 6.25"
```

Booleans

Python implements all of the usual operators for Boolean logic, but uses English words rather than symbols (&&, ||, etc.):

```
t = True
f = False
print(type(t)) # Prints "<class 'bool'>"
print(t and f) # Logical AND; prints "False"
print(t or f)  # Logical OR; prints "True"
print(not t)   # Logical NOT; prints "False"
print(t != f)  # Logical XOR; prints "True"
```

Strings

Python has great support for strings:

```
hello = 'hello'    # String literals can use single quotes
world = "world"    # or double quotes; it does not matter.
print(hello)       # Prints "hello"
print(len(hello))  # String length; prints "5"
hw = hello + ' ' + world # String concatenation
print(hw)          # prints "hello world"
```



```
hw12 = '%s %s %d' % (hello, world, 12) # sprintf style string formatting
print(hw12) # prints "hello world 12"
```

String objects have a bunch of useful methods; for example:

```
s = "hello"
print(s.capitalize()) # Capitalize a string; prints "Hello"
print(s.upper())      # Convert a string to uppercase; prints "HELLO"
print(s.rjust(7))     # Right-justify a string, padding with spaces;
prints "  hello"
print(s.center(7))    # Center a string, padding with spaces; prints "
hello "
print(s.replace('l', '(ell)')) # Replace all instances of one substring
with another;
                                # prints "he(ell)(ell)o"
print('  world '.strip()) # Strip leading and trailing whitespace;
prints "world"
```

You can find a list of all string methods in the documentation -
(<https://docs.python.org/3.8/library/stdtypes.html#string-methods>)

Lists

A list is the Python equivalent of an array, but is resizable and can contain elements of different types:

```
xs = [3, 1, 2] # Create a list
print(xs, xs[2]) # Prints "[3, 1, 2] 2"
print(xs[-1]) # Negative indices count from the end of the list;
prints "2"
xs[2] = 'foo' # Lists can contain elements of different types
print(xs) # Prints "[3, 1, 'foo']"
xs.append('bar') # Add a new element to the end of the list
print(xs) # Prints "[3, 1, 'foo', 'bar']"
x = xs.pop() # Remove and return the last element of the list
print(x, xs) # Prints "bar [3, 1, 'foo']"
```

More information is here <https://docs.python.org/3.8/tutorial/datastructures.html#more-on-lists>.

In addition to accessing list elements one at a time, Python provides concise syntax to access sublists; this is known as slicing:

```
nums = list(range(5)) # range is a built-in function that creates a
list of integers
print(nums) # Prints "[0, 1, 2, 3, 4]"
print(nums[2:4]) # Get a slice from index 2 to 4 (exclusive);
prints "[2, 3]"
print(nums[2:]) # Get a slice from index 2 to the end; prints
"[2, 3, 4]"
print(nums[:2]) # Get a slice from the start to index 2
(exclusive); prints "[0, 1]"
```



```
print(nums[:])           # Get a slice of the whole list; prints "[0, 1, 2, 3, 4]"
print(nums[:-1])         # Slice indices can be negative; prints "[0, 1, 2, 3]"
nums[2:4] = [8, 9]       # Assign a new sublist to a slice
print(nums)              # Prints "[0, 1, 8, 9, 4]"
```

We will see slicing again in the context of numpy arrays.

You can loop over the elements of a list like this:

```
animals = ['cat', 'dog', 'monkey']
for animal in animals:
    print(animal)
# Prints "cat", "dog", "monkey", each on its own line.
```

If you want access to the index of each element within the body of a loop, use the built-in `enumerate` function:

```
animals = ['cat', 'dog', 'monkey']
for idx, animal in enumerate(animals):
    print('#%d: %s' % (idx + 1, animal))
# Prints "#1: cat", "#2: dog", "#3: monkey", each on its own line
```

When programming, frequently we want to transform one type of data into another. As a simple example, consider the following code that computes square numbers:

```
nums = [0, 1, 2, 3, 4]
squares = []
for x in nums:
    squares.append(x ** 2)
print(squares)  # Prints [0, 1, 4, 9, 16]
```

You can make this code simpler using a list comprehension:

```
nums = [0, 1, 2, 3, 4]
squares = [x ** 2 for x in nums]
print(squares)  # Prints [0, 1, 4, 9, 16]
```

List comprehensions can also contain conditions:

```
nums = [0, 1, 2, 3, 4]
even_squares = [x ** 2 for x in nums if x % 2 == 0]
print(even_squares)  # Prints "[0, 4, 16]"
```

Dictionaries

A dictionary stores (key, value) pairs, similar to a Map in Java or an object in Javascript. You can use it like this:



```
d = {'cat': 'cute', 'dog': 'furry'} # Create a new dictionary with some
data
print(d['cat'])                    # Get an entry from a dictionary; prints "cute"
print('cat' in d)                  # Check if a dictionary has a given key; prints
"True"
d['fish'] = 'wet'                  # Set an entry in a dictionary
print(d['fish'])                   # Prints "wet"
# print(d['monkey'])               # KeyError: 'monkey' not a key of d
print(d.get('monkey', 'N/A'))     # Get an element with a default; prints
"N/A"
print(d.get('fish', 'N/A'))       # Get an element with a default; prints
"wet"
del d['fish']                      # Remove an element from a dictionary
print(d.get('fish', 'N/A'))       # "fish" is no longer a key; prints "N/A"
```

More information is here <https://docs.python.org/3.8/library/stdtypes.html#dict>.

It is easy to iterate over the keys in a dictionary:

```
d = {'person': 2, 'cat': 4, 'spider': 8}
for animal in d:
    legs = d[animal]
    print('A %s has %d legs' % (animal, legs))
# Prints "A person has 2 legs", "A cat has 4 legs", "A spider has 8 legs"
```

If you want access to keys and their corresponding values, use the `items` method:

```
d = {'person': 2, 'cat': 4, 'spider': 8}
for animal, legs in d.items():
    print('A %s has %d legs' % (animal, legs))
# Prints "A person has 2 legs", "A cat has 4 legs", "A spider has 8 legs"
```

These are similar to list comprehensions, but allow you to easily construct dictionaries. For example:

```
nums = [0, 1, 2, 3, 4]
even_num_to_square = {x: x ** 2 for x in nums if x % 2 == 0}
print(even_num_to_square) # Prints "{0: 0, 2: 4, 4: 16}"
```

Sets

A set is an unordered collection of distinct elements. As a simple example, consider the following:

```
animals = {'cat', 'dog'}
print('cat' in animals) # Check if an element is in a set; prints
"True"
print('fish' in animals) # prints "False"
animals.add('fish')       # Add an element to a set
print('fish' in animals)  # Prints "True"
print(len(animals))       # Number of elements in a set; prints "3"
```



```
animals.add('cat')           # Adding an element that is already in the set
                              # does nothing
print(len(animals))          # Prints "3"
animals.remove('cat')         # Remove an element from a set
print(len(animals))          # Prints "2"
```

More information is here <https://docs.python.org/3.8/library/stdtypes.html#set>

Iterating over a set has the same syntax as iterating over a list; however since sets are unordered, you cannot make assumptions about the order in which you visit the elements of the set:

```
animals = {'cat', 'dog', 'fish'}
for idx, animal in enumerate(animals):
    print('#%d: %s' % (idx + 1, animal))
# Prints "#1: fish", "#2: dog", "#3: cat"
```

Like lists and dictionaries, we can easily construct sets using set comprehensions:

```
from math import sqrt
nums = {int(sqrt(x)) for x in range(30)}
print(nums) # Prints "{0, 1, 2, 3, 4, 5}"
```

Tuples

A tuple is an (immutable) ordered list of values. A tuple is in many ways similar to a list; one of the most important differences is that tuples can be used as keys in dictionaries and as elements of sets, while lists cannot. Here is a trivial example:

```
d = {(x, x + 1): x for x in range(10)} # Create a dictionary with tuple
keys
t = (5, 6)                             # Create a tuple
print(type(t))                         # Prints "<class 'tuple'>"
print(d[t])                           # Prints "5"
print(d[(1, 2)])                      # Prints "1"
```

<https://docs.python.org/3.8/tutorial/datastructures.html#tuples-and-sequences> has more information about tuples.

Numpy

Numpy is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.

Arrays

A numpy array is a grid of values, all of the same type, and is indexed by a tuple of non negative integers. The number of dimensions is the rank of the array; the shape of an array is a tuple of integers giving the size of the array along each dimension.

We can initialize numpy arrays from nested Python lists, and access elements using square brackets:

```
import numpy as np
```



```
a = np.array([1, 2, 3])    # Create a rank 1 array
print(type(a))            # Prints "<class 'numpy.ndarray'>"
print(a.shape)            # Prints "(3,)"
print(a[0], a[1], a[2])   # Prints "1 2 3"
a[0] = 5                  # Change an element of the array
print(a)                  # Prints "[5, 2, 3]"

b = np.array([[1,2,3],[4,5,6]]) # Create a rank 2 array
print(b.shape)            # Prints "(2, 3)"
print(b[0, 0], b[0, 1], b[1, 0]) # Prints "1 2 4"
```

Numpy also provides many functions to create arrays:

```
import numpy as np

a = np.zeros((2,2))        # Create an array of all zeros
print(a)                  # Prints "[[ 0.  0.]
                          #           [ 0.  0.]]"

b = np.ones((1,2))         # Create an array of all ones
print(b)                  # Prints "[[ 1.  1.]]"

c = np.full((2,2), 7)      # Create a constant array
print(c)                  # Prints "[[ 7.  7.]
                          #           [ 7.  7.]]"

d = np.eye(2)              # Create a 2x2 identity matrix
print(d)                  # Prints "[[ 1.  0.]
                          #           [ 0.  1.]]"

e = np.random.random((2,2)) # Create an array filled with random values
print(e)                  # Might print "[[ 0.91940167  0.08143941]
                          #           [ 0.68744134  0.87236687]]"
```

More information here - <https://numpy.org/doc/stable/user/basics.creation.html#arrays-creation>

Array indexing

Numpy offers several ways to index into arrays. Similar to Python lists, numpy arrays can be sliced. Since arrays may be multidimensional, you must specify a slice for each dimension of the array:

```
import numpy as np

# Create the following rank 2 array with shape (3, 4)
# [[ 1  2  3  4]
#  [ 5  6  7  8]
#  [ 9 10 11 12]]
a = np.array([[1,2,3,4], [5,6,7,8], [9,10,11,12]])
```



```
# Use slicing to pull out the subarray consisting of the first 2 rows
# and columns 1 and 2; b is the following array of shape (2, 2):
# [[2 3]
#  [6 7]]
b = a[:2, 1:3]

# A slice of an array is a view into the same data, so modifying it
# will modify the original array.
print(a[0, 1])    # Prints "2"
b[0, 0] = 77      # b[0, 0] is the same piece of data as a[0, 1]
print(a[0, 1])    # Prints "77"
```

You can also mix integer indexing with slice indexing. However, doing so will yield an array of lower rank than the original array.

```
import numpy as np

# Create the following rank 2 array with shape (3, 4)
# [[ 1  2  3  4]
#  [ 5  6  7  8]
#  [ 9 10 11 12]]
a = np.array([[1,2,3,4], [5,6,7,8], [9,10,11,12]])

# Two ways of accessing the data in the middle row of the array.
# Mixing integer indexing with slices yields an array of lower rank,
# while using only slices yields an array of the same rank as the
# original array:
row_r1 = a[1, :]    # Rank 1 view of the second row of a
row_r2 = a[1:2, :]  # Rank 2 view of the second row of a
print(row_r1, row_r1.shape)    # Prints "[5 6 7 8] (4,)"
print(row_r2, row_r2.shape)    # Prints "[[5 6 7 8]] (1, 4)"

# We can make the same distinction when accessing columns of an array:
col_r1 = a[:, 1]
col_r2 = a[:, 1:2]
print(col_r1, col_r1.shape)    # Prints "[ 2  6 10] (3,)"
print(col_r2, col_r2.shape)    # Prints "[[ 2]
                                #          [ 6]
                                #          [10]] (3, 1)"
```

When you index into numpy arrays using slicing, the resulting array view will always be a subarray of the original array. In contrast, integer array indexing allows you to construct arbitrary arrays using the data from another array. Here is an example:

```
import numpy as np

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




```
# An example of integer array indexing.
# The returned array will have shape (3,) and
print(a[[0, 1, 2], [0, 1, 0]]) # Prints "[1 4 5]"

# The above example of integer array indexing is equivalent to this:
print(np.array([a[0, 0], a[1, 1], a[2, 0]])) # Prints "[1 4 5]"

# When using integer array indexing, you can reuse the same
# element from the source array:
print(a[[0, 0], [1, 1]]) # Prints "[2 2]"

# Equivalent to the previous integer array indexing example
print(np.array([a[0, 1], a[0, 1]])) # Prints "[2 2]"
```

One useful trick with integer array indexing is selecting or mutating one element from each row of a matrix:

```
import numpy as np

# Create a new array from which we will select elements
a = np.array([[1,2,3], [4,5,6], [7,8,9], [10, 11, 12]])

print(a) # prints "array([[ 1,  2,  3],
#                [ 4,  5,  6],
#                [ 7,  8,  9],
#                [10, 11, 12]])"

# Create an array of indices
b = np.array([0, 2, 0, 1])

# Select one element from each row of a using the indices in b
print(a[np.arange(4), b]) # Prints "[ 1  6  7 11]"

# Mutate one element from each row of a using the indices in b
a[np.arange(4), b] += 10

print(a) # prints "array([[11,  2,  3],
#                [ 4,  5, 16],
#                [17,  8,  9],
#                [10, 21, 12]])"
```

Boolean array indexing lets you pick out arbitrary elements of an array. Frequently this type of indexing is used to select the elements of an array that satisfy some condition. Here is an example:

```
import numpy as np

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



```
bool_idx = (a > 2)    # Find the elements of a that are bigger than 2;
                     # this returns a numpy array of Booleans of the same
                     # shape as a, where each slot of bool_idx tells
                     # whether that element of a is > 2.

print(bool_idx)      # Prints "[False False]
                     #           [ True  True]
                     #           [ True  True]]"
```

```
# We use boolean array indexing to construct a rank 1 array
# consisting of the elements of a corresponding to the True values
# of bool_idx
print(a[bool_idx])   # Prints "[3 4 5 6]"

# We can do all of the above in a single concise statement:
print(a[a > 2])      # Prints "[3 4 5 6]"
```

For brevity we have left out a lot of details about numpy array indexing; if you want to know more you should read <https://numpy.org/doc/stable/reference/arrays.indexing.html>

Every numpy array is a grid of elements of the same type. Numpy provides a large set of numeric datatypes that you can use to construct arrays. Numpy tries to guess a datatype when you create an array, but functions that construct arrays usually also include an optional argument to explicitly specify the datatype. Here is an example:

```
import numpy as np

x = np.array([1, 2])    # Let numpy choose the datatype
print(x.dtype)         # Prints "int64"

x = np.array([1.0, 2.0]) # Let numpy choose the datatype
print(x.dtype)         # Prints "float64"

x = np.array([1, 2], dtype=np.int64) # Force a particular datatype
print(x.dtype)         # Prints "int64"
```

More information is here <https://numpy.org/doc/stable/reference/arrays.dtypes.html>

Array math

Basic mathematical functions operate elementwise on arrays, and are available both as operator overloads and as functions in the numpy module:

```
import numpy as np

x = np.array([[1,2],[3,4]], dtype=np.float64)
y = np.array([[5,6],[7,8]], dtype=np.float64)
```



```
# Elementwise sum; both produce the array
# [[ 6.0  8.0]
#  [10.0 12.0]]
print(x + y)
print(np.add(x, y))

# Elementwise difference; both produce the array
# [[-4.0 -4.0]
#  [-4.0 -4.0]]
print(x - y)
print(np.subtract(x, y))

# Elementwise product; both produce the array
# [[ 5.0 12.0]
#  [21.0 32.0]]
print(x * y)
print(np.multiply(x, y))

# Elementwise division; both produce the array
# [[ 0.2          0.33333333]
#  [ 0.42857143  0.5         ]]
print(x / y)
print(np.divide(x, y))

# Elementwise square root; produces the array
# [[ 1.          1.41421356]
#  [ 1.73205081  2.         ]]
print(np.sqrt(x))
```

Note that `*` is elementwise multiplication, not matrix multiplication. We instead use the `dot` function to compute inner products of vectors, to multiply a vector by a matrix, and to multiply matrices. `dot` is available both as a function in the numpy module and as an instance method of array objects:

```
import numpy as np

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

v = np.array([9,10])
w = np.array([11, 12])

# Inner product of vectors; both produce 219
print(v.dot(w))
print(np.dot(v, w))

# Matrix / vector product; both produce the rank 1 array [29 67]
print(x.dot(v))
```



```
print(np.dot(x, v))

# Matrix / matrix product; both produce the rank 2 array
# [[19 22]
#  [43 50]]
print(x.dot(y))
print(np.dot(x, y))
```

Numpy provides many useful functions for performing computations on arrays; one of the most useful is `sum`:

```
import numpy as np

x = np.array([[1,2],[3,4]])

print(np.sum(x)) # Compute sum of all elements; prints "10"
print(np.sum(x, axis=0)) # Compute sum of each column; prints "[4 6]"
print(np.sum(x, axis=1)) # Compute sum of each row; prints "[3 7]"
```

You can find the full list of mathematical functions provided by numpy - <https://numpy.org/doc/stable/reference/routines.math.html>

Apart from computing mathematical functions using arrays, we frequently need to reshape or otherwise manipulate data in arrays. The simplest example of this type of operation is transposing a matrix; to transpose a matrix, simply use the `T` attribute of an array object:

```
import numpy as np

x = np.array([[1,2], [3,4]])
print(x)      # Prints "[[1 2]
               #           [3 4]]"
print(x.T)    # Prints "[[1 3]
               #           [2 4]]"

# Note that taking the transpose of a rank 1 array does nothing:
v = np.array([1,2,3])
print(v)      # Prints "[1 2 3]"
print(v.T)    # Prints "[1 2 3]"
```

Numpy provides many more functions for manipulating arrays; you can see the full list <https://numpy.org/doc/stable/reference/routines.array-manipulation.html>

Broadcasting

Broadcasting is a powerful mechanism that allows numpy to work with arrays of different shapes when performing arithmetic operations. Frequently we have a smaller array and a larger array, and we want to use the smaller array multiple times to perform some operation on the larger array.

For example, suppose that we want to add a constant vector to each row of a matrix. We could do it like this:



```
import numpy as np

# We will add the vector v to each row of the matrix x,
# storing the result in the matrix y
x = np.array([[1,2,3], [4,5,6], [7,8,9], [10, 11, 12]])
v = np.array([1, 0, 1])
y = np.empty_like(x)    # Create an empty matrix with the same shape as x

# Add the vector v to each row of the matrix x with an explicit loop
for i in range(4):
    y[i, :] = x[i, :] + v

# Now y is the following
# [[ 2  2  4]
#  [ 5  5  7]
#  [ 8  8 10]
#  [11 11 13]]
print(y)
```

This works; however when the matrix *x* is very large, computing an explicit loop in Python could be slow. Note that adding the vector *v* to each row of the matrix *x* is equivalent to forming a matrix *vv* by stacking multiple copies of *v* vertically, then performing elementwise summation of *x* and *vv*. We could implement this approach like this:

```
import numpy as np

# We will add the vector v to each row of the matrix x,
# storing the result in the matrix y
x = np.array([[1,2,3], [4,5,6], [7,8,9], [10, 11, 12]])
v = np.array([1, 0, 1])
vv = np.tile(v, (4, 1))    # Stack 4 copies of v on top of each other
print(vv)                  # Prints "[[ 1  0  1]
                            #          [ 1  0  1]
                            #          [ 1  0  1]
                            #          [ 1  0  1]]"

y = x + vv    # Add x and vv elementwise
print(y)      # Prints "[[ 2  2  4]
                #          [ 5  5  7]
                #          [ 8  8 10]
                #          [11 11 13]]"
```

Numpy broadcasting allows us to perform this computation without actually creating multiple copies of *v*. Consider this version, using broadcasting:

```
import numpy as np

# We will add the vector v to each row of the matrix x,
```



```
# storing the result in the matrix y
x = np.array([[1,2,3], [4,5,6], [7,8,9], [10, 11, 12]])
v = np.array([1, 0, 1])
y = x + v # Add v to each row of x using broadcasting
print(y) # Prints "[[ 2  2  4]
          #          [ 5  5  7]
          #          [ 8  8 10]
          #          [11 11 13]]"
```

The line `y = x + v` works even though `x` has shape (4, 3) and `v` has shape (3,) due to broadcasting; this line works as if `v` actually had shape (4, 3), where each row was a copy of `v`, and the sum was performed elementwise.

Broadcasting two arrays together follows these rules:

1. If the arrays do not have the same rank, prepend the shape of the lower rank array with 1s until both shapes have the same length.
2. The two arrays are said to be compatible in a dimension if they have the same size in the dimension, or if one of the arrays has size 1 in that dimension.
3. The arrays can be broadcast together if they are compatible in all dimensions.
4. After broadcasting, each array behaves as if it had shape equal to the elementwise maximum of shapes of the two input arrays.
5. In any dimension where one array had size 1 and the other array had size greater than 1, the first array behaves as if it were copied along that dimension

Here are some applications of broadcasting:

```
import numpy as np

# Compute outer product of vectors
v = np.array([1,2,3]) # v has shape (3,)
w = np.array([4,5])   # w has shape (2,)
# To compute an outer product, we first reshape v to be a column
# vector of shape (3, 1); we can then broadcast it against w to yield
# an output of shape (3, 2), which is the outer product of v and w:
# [[ 4  5]
#  [ 8 10]
#  [12 15]]
print(np.reshape(v, (3, 1)) * w)

# Add a vector to each row of a matrix
x = np.array([[1,2,3], [4,5,6]])
# x has shape (2, 3) and v has shape (3,) so they broadcast to (2, 3),
# giving the following matrix:
# [[2 4 6]
#  [5 7 9]]
print(x + v)
```



```
# Add a vector to each column of a matrix
# x has shape (2, 3) and w has shape (2,).
# If we transpose x then it has shape (3, 2) and can be broadcast
# against w to yield a result of shape (3, 2); transposing this result
# yields the final result of shape (2, 3) which is the matrix x with
# the vector w added to each column. Gives the following matrix:
# [[ 5  6  7]
#  [ 9 10 11]]
print((x.T + w).T)

# Another solution is to reshape w to be a column vector of shape (2, 1);
# we can then broadcast it directly against x to produce the same
# output.
print(x + np.reshape(w, (2, 1)))

# Multiply a matrix by a constant:
# x has shape (2, 3). Numpy treats scalars as arrays of shape ();
# these can be broadcast together to shape (2, 3), producing the
# following array:
# [[ 2  4  6]
#  [ 8 10 12]]
print(x * 2)
```

Broadcasting typically makes your code more concise and faster, so you should strive to use it where possible.

SciPy

Numpy provides a high-performance multidimensional array and basic tools to compute with and manipulate these arrays. SciPy builds on this, and provides a large number of functions that operate on numpy arrays and are useful for different types of scientific and engineering applications.

The best way to get familiar with SciPy is to browse the documentation - <https://docs.scipy.org/doc/scipy/reference/index.html>. We will highlight some parts of SciPy that you might find useful for this class.

Distance between points

SciPy defines some useful functions for computing distances between sets of points.

The function `scipy.spatial.distance.pdist` computes the distance between all pairs of points in a given set:

```
import numpy as np
from scipy.spatial.distance import pdist, squareform

# Create the following array where each row is a point in 2D space:
# [[0 1]
#  [1 0]
#  [2 0]]
x = np.array([[0, 1], [1, 0], [2, 0]])
print(x)
```



```
# Compute the Euclidean distance between all rows of x.
# d[i, j] is the Euclidean distance between x[i, :] and x[j, :],
# and d is the following array:
# [[ 0.          1.41421356  2.23606798]
#   [ 1.41421356  0.          1.          ]
#   [ 2.23606798  1.          0.          ]]
d = squareform(pdist(x, 'euclidean'))
print(d)
```

A similar function (`scipy.spatial.distance.cdist`) computes the distance between all pairs across two sets of points; you can read about it in the documentation -

<https://docs.scipy.org/doc/scipy/reference/generated/scipy.spatial.distance.cdist.html>

Matplotlib

Matplotlib is a plotting library. In this section give a brief introduction to the `matplotlib.pyplot` module, which provides a plotting system similar to that of MATLAB.

The most important function in `matplotlib` is `plot`, which allows you to plot 2D data. Here is a simple example:

```
import numpy as np
import matplotlib.pyplot as plt

# Compute the x and y coordinates for points on a sine curve
x = np.arange(0, 3 * np.pi, 0.1)
y = np.sin(x)

# Plot the points using matplotlib
plt.plot(x, y)
plt.show() # You must call plt.show() to make graphics appear.
```

With just a little bit of extra work we can easily plot multiple lines at once, and add a title, legend, and axis labels:

```
import numpy as np
import matplotlib.pyplot as plt

# Compute the x and y coordinates for points on sine and cosine curves
x = np.arange(0, 3 * np.pi, 0.1)
y_sin = np.sin(x)
y_cos = np.cos(x)

# Plot the points using matplotlib
plt.plot(x, y_sin)
plt.plot(x, y_cos)
plt.xlabel('x axis label')
plt.ylabel('y axis label')
plt.title('Sine and Cosine')
```




```
plt.legend(['Sine', 'Cosine'])  
plt.show()
```

You can read more - https://matplotlib.org/2.0.2/api/pyplot_api.html

You can plot different things in the same figure using the subplot function. Here is an example:

```
import numpy as np  
import matplotlib.pyplot as plt  
  
# Compute the x and y coordinates for points on sine and cosine curves  
x = np.arange(0, 3 * np.pi, 0.1)  
y_sin = np.sin(x)  
y_cos = np.cos(x)  
  
# Set up a subplot grid that has height 2 and width 1,  
# and set the first such subplot as active.  
plt.subplot(2, 1, 1)  
  
# Make the first plot  
plt.plot(x, y_sin)  
plt.title('Sine')  
  
# Set the second subplot as active, and make the second plot.  
plt.subplot(2, 1, 2)  
plt.plot(x, y_cos)  
plt.title('Cosine')  
  
# Show the figure.  
plt.show()
```

Seaborn

We are going to use the iris flowers dataset. This dataset is famous because it is used as the “hello world” dataset in machine learning and statistics by pretty much everyone. The dataset contains 150 observations of iris flowers. There are four columns of measurements of the flowers in centimeters. The fifth column is the species of the flower observed. All observed flowers belong to one of three species. We can load the data directly from the UCI Machine Learning repository. We are using pandas to load the data.

```
# summarize the data  
from pandas import read_csv  
# Load dataset  
url =  
"https://raw.githubusercontent.com/jbrownlee/Datasets/master/iris.csv"  
names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width',  
         'class']  
dataset = read_csv(url, names=names)  
# shape
```



```
print(dataset.shape)
# head
print(dataset.head(20))
# descriptions
print(dataset.describe())
# class distribution
print(dataset.groupby('class').size())
```

We start with checking data distribution:

```
import seaborn as sns
import matplotlib.pyplot as plt

sns.jointplot(x='sepal-length', y='sepal-width', data=dataset, height=5)
sns.FacetGrid(dataset, hue='class', height=5).map(plt.scatter, 'sepal -
length', 'sepal-width').add_legend()
plt.show() # Display the plot
```

Scikit-learn

Next create some models of the data and estimate their accuracy on unseen data:

- Separate out a validation dataset.
- Set-up the test harness to use 10-fold cross validation.
- Build two different models (using scikit-learn algorithms) to predict species from flower measurements.
- Select the best model.

```
# Compare algorithms
from pandas import read_csv
from matplotlib import pyplot
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import StratifiedKFold
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier

# Load dataset
url =
"https://raw.githubusercontent.com/jbrownlee/Datasets/master/iris.csv"
names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width',
'class']
dataset = read_csv(url, names=names)

# Split-out validation dataset
array = dataset.values
```



```
X = array[:,0:4]
y = array[:,4]
X_train, X_validation, Y_train, Y_validation = train_test_split(X, y,
test_size=0.20, random_state=1, shuffle=True)

# Spot Check Algorithms
models = []
models.append(('LR', LogisticRegression(solver='liblinear',
multi_class='ovr'))
models.append(('CART', DecisionTreeClassifier()))

# Evaluate each model in turn
results = []
names = []
for name, model in models:
    kfold = StratifiedKFold(n_splits=10, random_state=1, shuffle=True)
    cv_results = cross_val_score(model, X_train, Y_train, cv=kfold,
scoring='accuracy')
    results.append(cv_results)
    names.append(name)
    print('%s: %f (%f)' % (name, cv_results.mean(), cv_results.std()))

# Compare Algorithms
pyplot.boxplot(results, labels=names)
pyplot.title('Algorithm Comparison')
pyplot.show()
```