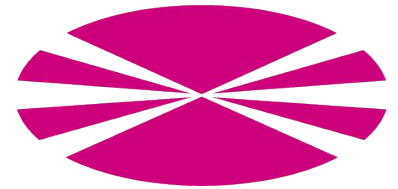


T3: Introduction to Data Analytics

Data Analytics with HPC
Master in High Performance Computing, Spring 2025



Contents

- 6/2 - Introduction to Python
 - Execution environments: iPython, Jupyter.
 - Basics.
 - Operators and types. Instructions. Functions. Modules. Classes. Exceptions.
 - Integrated tools.
 - NumPy.
 - Matplotlib.

Contents

- 13/2 - Data Analytics: Pandas
 - Intro to Pandas.
 - Data I/O.
 - Data wrangling.
 - Visualization.
 - Aggregation.
 - Time series.

Contents

- 20/2 - High level tools
 - Visualization: Seaborn.
 - Machine learning: Scikit-Learn
 - Out-of-core computing: Dask

Contents

- 27/2 - Projects
 - Time for developing your data analysis project.

Introduction to Python

Contents

- Execution environments: iPython, Jupyter.
- Basics.
- Operators and types. Instructions. Functions. Modules. Classes. Exceptions.
- Integrated tools.
- NumPy.
- Matplotlib.

Python



- General-purpose language, interpreted, object-oriented.
- Portable: an interpreter exists for almost any OS.
- Open-source reference implementation (CPython), but alternative ones exist as well (e.g., iPython, Intel Python).
- Powerful and versatile:
 - Dynamic typing.
 - Complex “basic” types: lists, dictionaries, strings, ...
 - Built-in tools: powerful list operations, strings, arrays, ...
 - Libraries: math, statistics, parallelization, ...
 - Automatic memory management.

Python

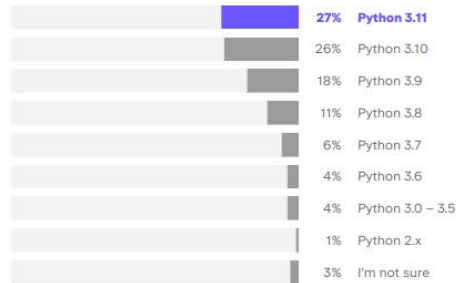


- Popular with the scientific and engineering communities.
- Free and general purpose alternative to Matlab/R.
- “Slice” notation (Matlab-like).
- Multiple extensions: NumPy, Matplotlib, SciPy, pytables, ...
- Simple integration: C/C++/Fortran, web services, ...
- Web development: Django, Flask, Pyramid, ...
- Scientific computing environments: iPython, Anaconda, ...

Python: “state of the art”

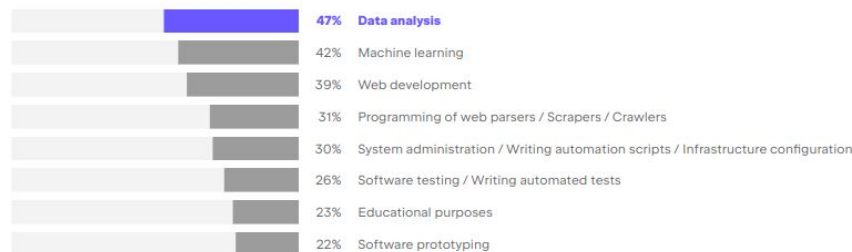


Which version of Python do you use the most?



One in four respondents use the latest version of Python, released in October 2022. It took Python 3.11 approximately nine months to reach 27% adoption.

What do you use Python for?



Show more

The use cases of Python have remained stable year over year, with the most popular areas of usage being data analysis, web development, and machine learning.

Source: <https://www.jetbrains.com/lp/devecosystem-2023/python/>

iPython



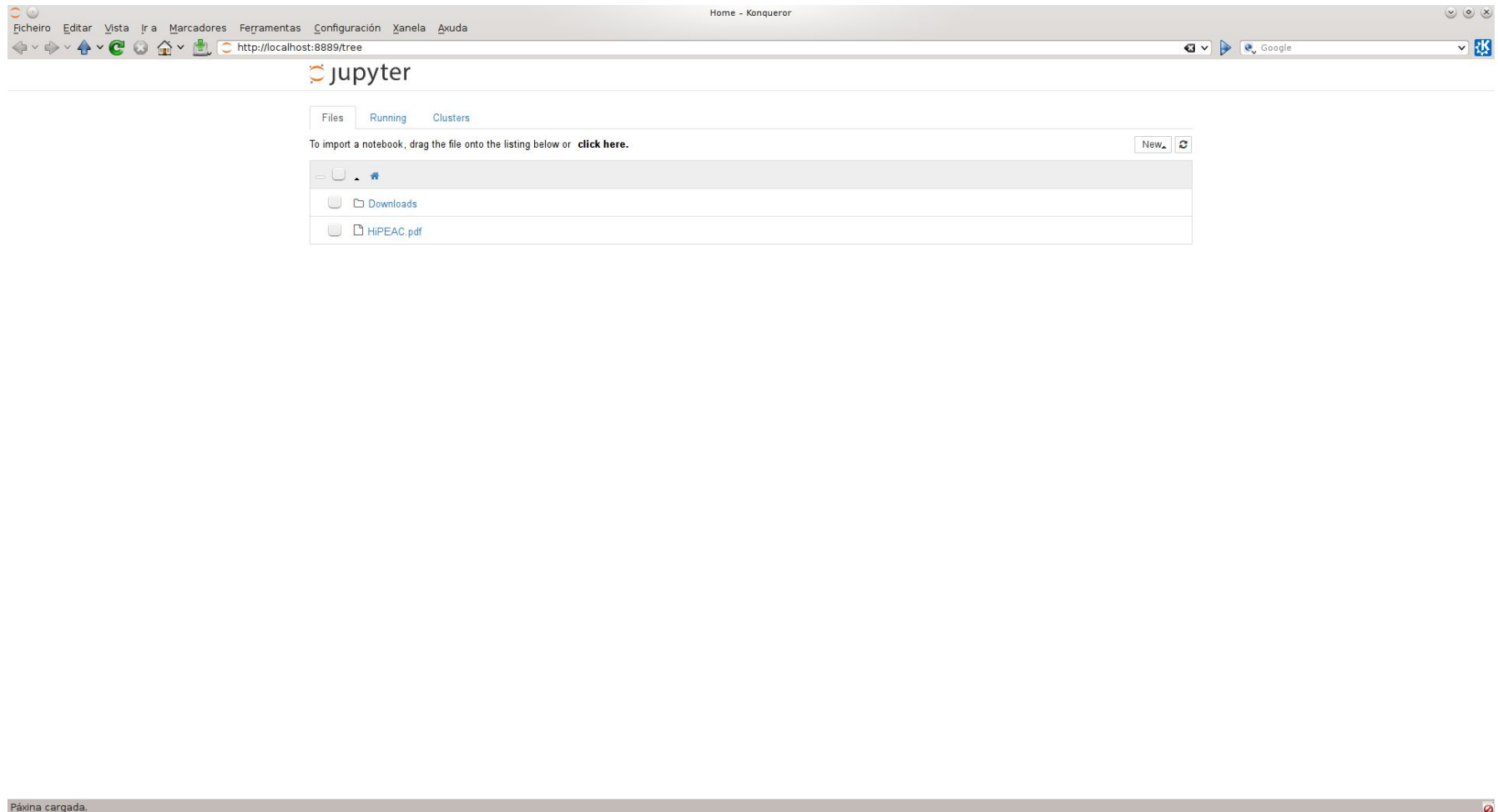
- Interactive Python interpreter.
- Does not provide analytics or computing tools, but greatly eases up the use of Python and the development process.
- Includes GUIs for easy data visualization.
- Allows to work remotely through a web browser, as well as storing work sessions in an HTML format to allow sharing code, data, and examples.

Jupyter Notebooks



- Environment focused on exporting HTML or being used interactively.
- Uses a JSON-based format to allow sharing code and outputs.
- Built as a web server to which clients are connected using a browser.
- Allows to remotely use a Python environment via web.

Jupyter Notebooks



Jupyter Notebooks



Untitled - Chromium

localhost:8888/notebooks/Untitled.ipynb?kernel_name=python2

Apps Euribor en el BOE Descargas Piano Cocina Fotografía Journals TUTORIAL - Recept Android Viajes Venues Essential Albums EURAXESS - Europ BIGMAT PROMACO PSU12 Trabajo Other bookmarks

jupyter Untitled Last Checkpoint: a minute ago (unsaved changes)

File Edit View Insert Cell Kernel Help Python 2

Cell Toolbar: None

```
In [1]: %pylab inline
        Populating the interactive namespace from numpy and matplotlib

In [2]: img=plt.imread( '/tmp/Mandelbrot-Javascript.png' )

In [3]: imshow(img)
Out[3]: <matplotlib.image.AxesImage at 0x7f847692b410>
```

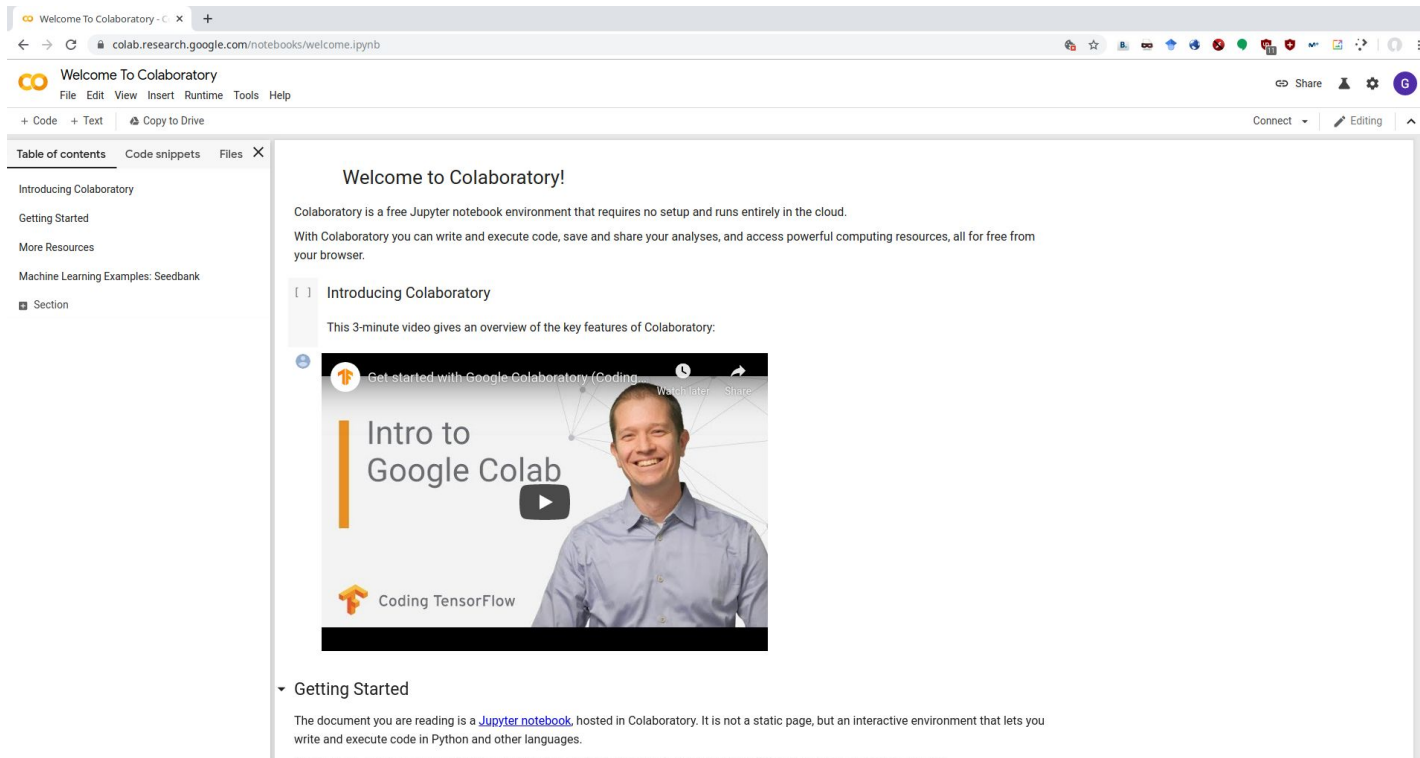
A visualization of the Mandelbrot set, showing the complex plane with the set's boundary highlighted in red and black. The axes are labeled from 0 to 600 on the y-axis and 0 to 800 on the x-axis.

In []:

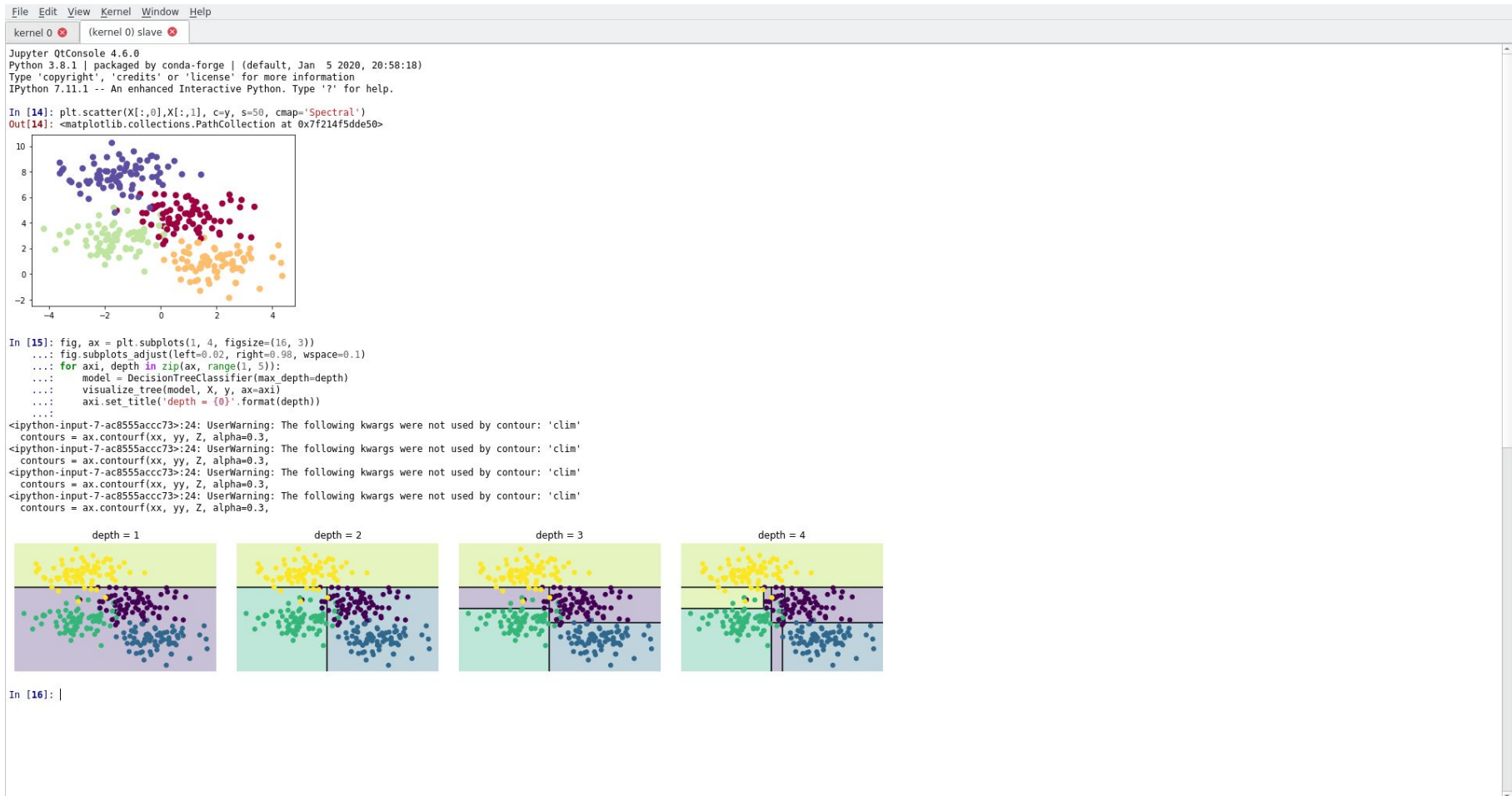
Google Colaboratory



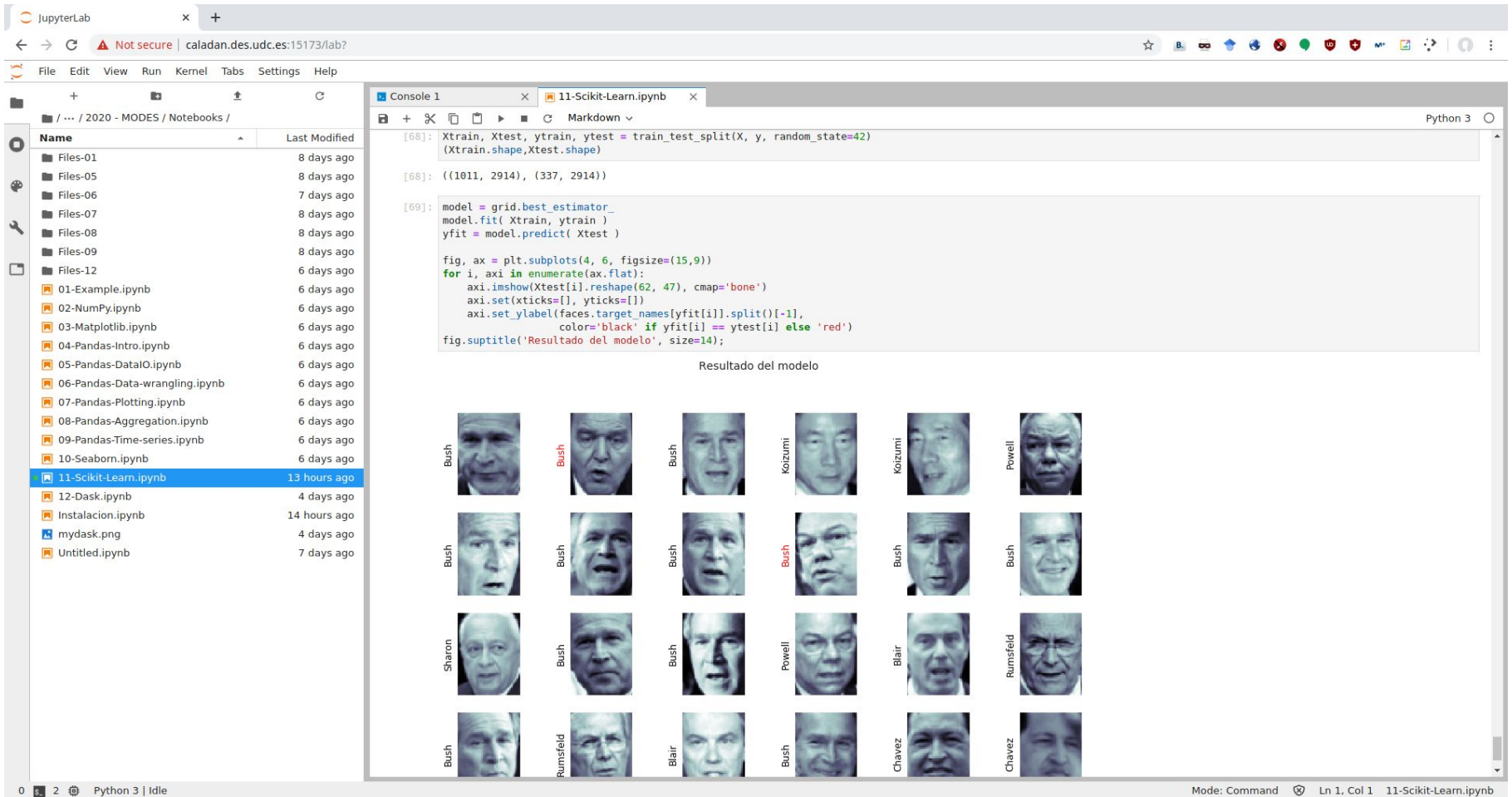
- Free notebook environment on the cloud.
- <https://colab.research.google.com>



Jupyter QtConsole



JupyterLab: Next Generation Notebook Interface



The screenshot displays the JupyterLab web interface. On the left is a file browser showing a directory structure with files like 'Files-01' through 'Files-12' and various '.ipynb' files. The main area is divided into two panes. The top pane shows a code editor with Python code for training a model and plotting results. The bottom pane shows the output of the code, which is a 4x6 grid of face images. Each image is labeled with a name, and some labels are highlighted in red. The console at the bottom shows the execution of the code and the resulting output.

```
[68]: Xtrain, Xtest, ytrain, ytest = train_test_split(X, y, random_state=42)
      (Xtrain.shape, Xtest.shape)

[68]: ((1011, 2914), (337, 2914))

[69]: model = grid.best_estimator_
      model.fit(Xtrain, ytrain)
      yfit = model.predict(Xtest)

      fig, ax = plt.subplots(4, 6, figsize=(15,9))
      for i, axi in enumerate(ax.flat):
          axi.imshow(Xtest[i].reshape(62, 47), cmap='bone')
          axi.set_xticks=[], yticks=[]
          axi.set_ylabel(faces.target_names[yfit[i]].split()[-1],
                        color='black' if yfit[i] == ytest[i] else 'red')
      fig.suptitle('Resultado del modelo', size=14)
```

Resultado del modelo

Bush	Bush	Bush	Koizumi	Koizumi	Powell
Bush	Bush	Bush	Bush	Bush	Bush
Sharon	Bush	Bush	Powell	Blair	Rumsfeld
Bush	Rumsfeld	Blair	Bush	Chavez	Chavez

Python: Conceptual hierarchy



- Programs are made of **modules**.
- Modules contain **functions**.
- Functions contain **instructions**.
- Instructions contain **expressions**.
- Expressions create and manipulate **objects**.
- Objects belong to a **class**.

Python: Conceptual hierarchy



```
>>>>>>> file.py
```

```
def f():
```

```
    i = 0
```

```
    while i < 100:
```

```
        print( i )
```

```
        i += 1
```

Python: Objects

- **Everything** is an object:
 - integer numbers -> *int*
 - floating point numbers -> *float*
 - character strings -> *str*
 - lists -> *list*
 - dictionaries -> *dict*
- The type of a variable can be queried using *type()*.
- Objects are containers which aggregate variables (*attributes*) and functions (*methods*).
- Objects belong to a *class*, which is an abstract description (or scheme) of the common traits of a family of objects.

Python - Executing programs

1. Using the interactive interpreter:

```
>>> txt = "Luke, I am your father"
```

```
>>> print( txt )
```

```
Luke, I am your father
```

Python - Executing programs



2. Interpreting Python code:

```
----- > script.py
```

```
txt = "Luke, I am your father"
```

```
print( txt )
```

```
-----
```

```
$ python script.py
```

```
Luke, I am your father
```

Python - Executing programs



3. Shell script:

```
----- > script.py
```

```
#!/usr/bin/python
```

```
txt = "Luke, I am your father"
```

```
print( txt )
```

```
-----
```

```
$ ./script.py
```

```
Luke, I am your father
```

Python - Executing programs



4. Embedding Python into C:

```
----- > test.c

#include <Python.h>

...

Py_Initialize();

PyRun_SimpleString( "txt = 'Luke, I am your father'" );

PyRun_SimpleString( "print( txt )" );

...

-----

$ gcc -I /usr/include/python3.8 test.c -lpython3.8

$ ./a.out

Luke, I am your father
```

Python - Semantics & Syntax

Indentation, not brackets



```
if x < threshold:
    print( "x less than threshold" )
else:
    print( "x greater than threshold" )
```


Python - Semantics & Syntax

Everything is an object



- Including numbers, strings, functions, classes, modules, ...
- Each object has a type and its own data.
- This provides flexibility, allowing to deal with any element in a generic way.

Python - Semantics & Syntax

Comments



```
# This is a dramatic moment.  
  
txt = "Luke, I am your father"  
  
print( txt )
```

Python - Semantics & Syntax

Function calls

- C syntax:

```
result = f( x, y, z )
```

- Class methods:

```
result = obj.f( x, y, z )
```

- Functions may be passed positional or keyword arguments:

```
result = f( x, y, z, tol=0.01, method="fast" )
```

Python - Semantics & Syntax

Variables and passing by reference

- Assignment to a variable generates a new reference to an object.
- The original object remains the same!

```
>>> a = [1,2,3]
```

```
>>> b = a
```

```
>>> a.append(4)
```

```
>>> b
```

```
[1,2,3,4]
```

Python - Semantics & Syntax

Variables and passing by reference

Consequently:

- Passing large variables is efficient.
- Functions may permanently modify their parameters!
- Except for **immutable** types. E.g.,

```
>>> a = 10
```

```
>>> b = a
```

```
>>> a = 20
```

```
>>> print b
```

```
10
```

Python - Semantics & Syntax

Dynamic references, strong typing

References (variables) have no static type:

```
>>> a = 10
```

```
>>> type(a)
```

```
<type 'int'>
```

```
>>> a = "10"
```

```
>>> type(a)
```

```
<type 'str'>
```

Python - Semantics & Syntax

Dynamic references, strong typing

Python has strong typing, nevertheless:

```
>>> 5 + "5"
```

```
TypeError: unsupported operand type(s) for +: 'int' and 'str'
```

```
>>> "5" + 5
```

```
TypeError: cannot concatenate 'str' and 'int' objects
```

```
>>> 5 + int("5")
```

```
10
```

```
>>> "5" + str(5)
```

```
'55'
```

Each Python object has a specific type (class). Implicit casting is only performed for `int` to floating point conversions.

Python - Semantics & Syntax



Dynamic references, strong typing

The `isinstance` operator is used to check whether an object belongs to a particular class:

```
>>> a = 4.5
```

```
>>> isinstance(a, int)
```

```
False
```

```
>>> isinstance(a, float)
```

```
True
```


Python - Semantics & Syntax

Methods and Attributes



Python objects have methods and attributes:

```
>>> a = 4.5
```

```
>>> dir(a)
```

```
['__abs__', '__add__', '__class__', '__coerce__', '__delattr__', '__div__', '__divmod__',  
 '__doc__', '__eq__', '__float__', '__floordiv__', '__format__', '__ge__',  
 '__getattr__', '__getformat__', '__getnewargs__', '__gt__', '__hash__', '__init__',  
 '__int__', '__le__', '__long__', '__lt__', '__mod__', '__mul__', '__ne__', '__neg__',  
 '__new__', '__nonzero__', '__pos__', '__pow__', '__radd__', '__rdiv__', '__rdivmod__',  
 '__reduce__', '__reduce_ex__', '__repr__', '__rfloordiv__', '__rmod__', '__rmul__',  
 '__rpow__', '__rsub__', '__rtruediv__', '__setattr__', '__setformat__', '__sizeof__',  
 '__str__', '__sub__', '__subclasshook__', '__truediv__', '__trunc__', 'as_integer_ratio',  
 'conjugate', 'fromhex', 'hex', 'imag', 'is_integer', 'real']
```

Python - Semantics & Syntax

Methods and Attributes

```
>>> a.hex()
```

```
'0x1.2000000000000p+2'
```

```
>>> hasattr( a, "hex" )
```

```
True
```

```
>>> getattr( a, "hex" )
```

```
<built-in method hex of float object at 0x1ed86a0>
```

```
>>> getattr( a, "hex" )()
```

```
'0x1.2000000000000p+2'
```

Python - Semantics & Syntax

Importing modules

- A Python module is a `.py` file including definitions and/or code:

```
----- > modulol.py
```

```
PI = 3.14159
```

```
def f(x):
```

```
    return x + 2
```

```
def g(a,b):
```

```
    return a + b
```

```
-----
```

Python - Semantics & Syntax



Binary operators

<code>a + b</code>	Addition	<code>a - b</code>	Subtraction
<code>a * b</code>	Multiplication	<code>a / b</code>	Division (Python 3)
<code>a // b</code>	Integer division (Python 3)	<code>a ** b</code>	Exponentiation
<code>a & b</code>	Bitwise AND	<code>a b</code>	Bitwise OR
<code>a ^ b</code>	Bitwise XOR	<code>a == b</code>	Equality
<code>a != b</code>	Inequality	<code>a < b, a <= b</code>	Less/Less or equal than
<code>a > b, a >= b</code>	Greater/Greater or equal than	<code>a is b</code>	Identity: True if <code>a</code> and <code>b</code> are the same object

Python - Scalar types



<code>None</code>	"null" in Python. Singleton class.
<code>str</code>	String type. ASCII in Python 2.x, Unicode in Python 3.
<code>unicode</code>	Unicode string.
<code>float</code>	Floating point number, double precision (64 bits).
<code>bool</code>	Logical, <code>True</code> or <code>False</code> .
<code>int</code>	Signed integer. 32 or 64 bits depending on platform. Arbitrary precision in Python 3.
<code>long</code>	Signed integer, arbitrary precision. Does not exist in Python 3.

Python - Scalar types

Numerical types

- Python 2: `int` (32- or 64-bits), `float` (64 bits) and `long` (arbitrary precision).
- Python 3: `int` (arbitrary precision), `float` (64 bits).
- Conversion between `int` and `long` is transparent in Python 2.
- Scientific notation accepted:

```
>>> fval = 6.78e-5
```

Python - Scalar types

Numerical types

- Division of two integers in Python 2 yields an integer (truncated, not rounded). For floating point division of integers in Python 2:

```
>>> fval = a_int / float(b_int)
```

- Division of integers in Python 3 yields a floating point number. For integer division of integers in Python 3:

```
>>> ival = a_int // b_int
```

- Complex numbers are written using `j` for the imaginary part:

```
>>> cval = 1 + 2j
```

Python - Scalar types

Strings

- String literals are written between single or double quotes:

```
>>> a = 'one way to write a string'
```

```
>>> b = "a different way"
```

- To write strings with line breaks, use three single quotes:

```
>>> c = '''
```

```
This is a larger string
```

```
Which spans more than one line
```

```
'''
```


Python - Scalar types

Strings

- Strings are immutable objects: once they are created, they cannot be modified:

```
>>> a = 'this is a string'
```

```
>>> a[10] = 'f'
```

```
TypeError: 'str' object does not support item assignment
```

```
>>> b = a.replace( 'string', 'larger string')
```

```
>>> b
```

```
'this is a larger string'
```

Python - Scalar types

Strings

- Strings can be created from different types using `str()`:

```
>>> a = 5.6
```

```
>>> s = str(a)
```

```
>>> s
```

```
'5.6'
```

Python - Scalar types

Strings

- A string is just a **collection** of characters (so not really a scalar type):

```
>>> s = 'python'
```

```
>>> list(s)
```

```
[ 'p', 'y', 't', 'h', 'o', 'n' ]
```

```
>>> s[:3]
```

```
'pyt'
```

Python - Scalar types

Strings

- The backslash (\) is used as the scape character, to encode special characters such as line breaks (\n) or Unicode characters.
- A string which should be interpreted literally (raw) can be marked using `r'string'`:

```
>>> s = r'no\special\characters\in\this\string'
```

```
>>> s
```

```
'no\special\characters\in\this\string'
```

Python - Scalar types

Strings

- Operator + applies string concatenation:

```
>>> a = 'first half of string '
```

```
>>> b = 'and second half of string'
```

```
>>> a + b
```

```
'First half of string and second half of string'
```

Python - Scalar types

Strings

- To specify that a string is encoded using Unicode, `u'string'` is used.
- To format a string there's a C-like template syntax:

```
>>> template = u "%.2f %s are worth %d€"
```

```
>>> template % (1.18018, u'american dollars', 1)
```

```
u'1.18 american dollars are worth 1\u20ac'
```

```
>>> print template % (1.18018, u'american dollars', 1)
```

```
'1.18 american dollars are worth 1€'
```

Python - Scalar types

Booleans

- The two boolean values are written `True` and `False`.
- They can be combined using the keywords `and` and `or`.

```
>>> True and True
```

```
True
```

```
>>> False and True
```

```
True
```

Python - Scalar types

Booleans

- All basic predefined types in Python, as well as any class implementing the `__nonzero__()` method, can be interpreted as `True` or `False` in a conditional clause.

```
>>> a = [1,2,3]
```

```
>>> if a: print "True"
```

```
'True'
```

```
>>> b = []
```

```
>>> if not b: print "False"
```

```
'False'
```


Python - Scalar types

None

- `None` is the null value in Python.
- If a function does not explicitly return a value, it returns `None` implicitly.
- Can be used as a default value for optional parameters to functions:

```
def add_and_maybe_multiply( a, b, c = None ):
    result = a + b
    if c is not None: result *= c
    return result
```

- `None` is not a reserved word, but an instance (singleton) of `NoneType`.

Python - Scalar types

Dates and times

- The `datetime` provides the types `datetime`, `date` and `time`.

```
>>> from datetime import datetime, date, time
```

```
dt = datetime( 2015, 04, 24, 12, 25, 32 )
```

```
>>> dt.day
```

```
24
```

```
>>> dt.time()
```

```
datetime.time( 12, 25,  
32 )
```

Python - Scalar types

Dates and times

- The `strftime()` method formats a `datetime` object as a string:

```
>>> dt.strftime( '%d/%m/%Y %H:%M' )
```

```
'24/04/2015 12:25'
```

- A `datetime` object can be created from a string using `strptime`:

```
>>> datetime.strptime( '20150424', '%Y%m%d' )
```

```
datetime.datetime( 2015, 4, 24, 0, 0 )
```

Python - Scalar types

Dates and times

- The subtraction of two `datetime` objects returns a `datetime.timedelta`:

```
>>> dt2 = datetime( 2015, 4, 25, 12, 25, 32 )
```

```
>>> dt2 - dt
```

```
datetime.timedelta(1)
```

- Operating a `datetime` object with a `timedelta` returns a new `datetime`:

```
>>> from datetime import timedelta
```

```
>>> dt + timedelta(days=1)
```

```
>>> datetime.datetime(2015, 4, 25, 12, 25, 32)
```

Python - Flow control

`if, elif, else`

These blocks behave like in most other languages:

```
if x < 0:

    print 'Negative'

elif x == 0:

    print 'Zero'

elif 0 < x < 5:

    print 'Positive less than 5'

else:

    print 'Positive and greater or equal than 5'
```

Python - Flow control

for loops

- A `for` loop is used to iterate over a collection (like a list or tuple):

```
for val in collection:
```

```
    # Do something with val
```

- `continue` skips to the next iteration
- `break` exits the current loop

Python - Flow control

`range()` and `xrange()`

- `range()` returns a list of equally spaced integers:

```
>>> range(10)
```

```
[0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
```

- Takes as parameters the start, end, and step of the sequence:

```
>>> range(5, 20, 3)
```

```
[5, 8, 11, 14, 17]
```

- Generates values in the interval `[start, end)`.

Python - Flow control

`range()` and `xrange()`

- `range()` is useful to code for loops with C semantics:

```
>>> for x in range(10):
```

```
...     print 2*x
```

```
0
```

```
2
```

```
4
```

```
8
```

```
16
```

```
18
```


Python - Flow control

`range()` and `xrange()`

- `xrange()` accepts the same parameters as `range()`, but returns a generator instead of building the full list in memory.
- It features lazy evaluation.

```
>>> xrange(10)
```

```
xrange(10)
```

```
>>> for x in xrange(10):
```

```
...     print 2*x
```

```
0
```

```
2
```

```
...
```

- It is preferable to use `xrange()` when working with large ranges.

Python - Flow control

while loops

- Iterate **while** a condition is met:

```
x = 256
```

```
total = 0
```

```
while x > 0:
```

```
    if total > 500: break
```

```
    total += x
```

```
    x = x // 2
```

Python - Flow control

pass

- `pass` is a no-op instruction in Python.
- It is sometimes required because white spaces in Python delimit execution blocks:

```
if x < 0:

    print "Negativo"

elif x == 0: pass

else:

    print "Positivo"
```

Python - Flow control

Exception handling

- Functions and operations can raise exceptions:

```
>>> 5 / 0
```

```
ZeroDivisionError: integer division or modulo by zero
```

- Exceptions can be handled to solve runtime errors.
- This allows to dynamically manage some selected error types.

Python - Flow control

Exception handling

```
def float_division( a, b ):  
  
    x = NaN  
  
    try:  
  
        x = a / float(b)  
  
    except:  
  
        print "Exception during division"  
  
    return x
```

Python - Flow control

Exception handling

- We can explicitly list the exception types managed by a `except` block:

```
def float_division( a, b ):  
  
    x = NaN  
  
    try:  
  
        x = a / float(b)  
  
    except ZeroDivisionError:  
  
        print "Division by zero"  
  
    return x
```

Python - Flow control

Exception handling

- Multiple exception types can be managed by the same block:

```
def float_division( a, b ):  
  
    x = NaN  
  
    try:  
  
        x = a / float(b)  
  
    except ZeroDivisionError, TypeError:  
  
        print "Exception during division"  
  
    return x
```

Python - Flow control

Exception handling

- A `finally` block is executed regardless of whether the `try` was successful:

```
f = open( path, 'w' )

try:

    write_to_file( f )

finally:

    f.close()
```


Python - Flow control

Exception handling

- An `else` block will only be executed if the `try` was successful:

```
f = open( path, 'w' )

try:

    write_to_file( f )

except: print 'Failure'

else: print 'Success'

finally: f.close()
```

Python - Collections

Tuples

- Tuple: unidimensional sequence of fixed size containing Python objects.

```
>>> tup = (2, 3, 4)
```

```
>>> nested_tup = ( (1, 2), (2, 3, 4) )
```

```
>>> tup_from_list = tuple( [2,3,4] )
```

```
>>> tup_from_iter = tuple( xrange(5) )
```

```
>>> tup_from_string = tuple( 'string' )
```

Python - Collections

Tuples

- The elements in a tuple are accessed using the `[]` operator:

```
>>> tup_from_iter[3]
```

```
3
```

```
>>> tup_from_iter[1:3]
```

```
(1, 2)
```

```
>>> tup_from_string[0]
```

```
's'
```

Python - Collections

Tuples

- Tuples are immutable:

```
>>> tup_from_iter[3] = 5
```

```
TypeError: 'tuple' object does not support item assignment
```

- However, objects inside a tuple can be mutable:

```
>>> tup = ( (1,2), [3,4,5] )
```

```
>>> tup[1].append(6)
```

```
>>> tup
```

```
((1, 2), [3, 4, 5, 6])
```

Python - Collections

Tuples

- Tuples are concatenated using the + operator:

```
>>> tup_from_iter + tup_from_string
```

```
(0, 1, 2, 3, 4, 's', 't', 'r', 'i', 'n', 'g')
```

- The * operator is consistent with the addition/concatenation semantics:

```
>>> tup_from_iter*3
```

```
(0, 1, 2, 3, 4, 0, 1, 2, 3, 4, 0, 1, 2, 3, 4)
```

Python - Collections

Tuples

- Note that the objects contained in tuples are not copied, but referenced:

```
>>> tup = ( (1, 2), (3, 4, 5) )
```

```
>>> tup2 = tup * 3
```

```
>>> tup[1].append(6)
```

```
>>> tup2
```

```
((1, 2), [3, 4, 5, 6], (1, 2), [3, 4, 5, 6], (1, 2), [3, 4, 5, 6])
```

Python - Collections

Tuples

- If the right hand side of an assignment is a tuple, Python tries to unpack it:

```
>>> tup = (1,2,3)
```

```
>>> a, b, c = tup
```

```
>>> b
```

```
2
```

Python - Collections

Tuples

- Nested tuples can be explicitly unpacked:

```
>>> tup = ( 1, 2, (3, 4) )
```

```
>>> a, b, (c, d) = tup
```

```
>>> c
```

```
3
```

```
>>> a, b = b, a
```

```
>>> a
```

```
2
```

- A common use of tuples is to code functions that return multiple values.

Python - Collections

Lists

- Unlike tuples, lists have variable length and are mutable:

```
>>> list_1 = [2, 3, 7, None]
```

```
>>> tup = ( 'a', 'b', 'c' )
```

```
>>> list_2 = list( tup )
```

```
>>> list_2[1] = 'd'
```

```
>>> list_2
```

```
[ 'a', 'd', 'c' ]
```

Python - Collections

Lists

- Elements are added at the end of a list using `append()`:

```
>>> list_2.append( 5 )
```

```
>>> list_2
```

```
[ 'a', 'd', 'c', 5 ]
```

- Elements can also be added at a specific place using `insert()`:

```
>>> list_2.insert( 3, None )
```

```
>>> list_2
```

```
[ 'a', 'd', 'c', None, 5 ]
```

Python - Collections

Lists

- The reverse operation to `insert()` is `pop()`:

```
>>> list_2.pop()
```

```
5
```

```
>>> list_2
```

```
[ 'a', 'd', 'c', None ]
```

```
>>> list_2.pop(2)
```

```
'c'
```

```
>>> list_2
```

```
[ 'a', 'd', None ]
```

Python - Collections

Lists

- Elements can be deleted from the list using `remove()`:

```
>>> list_2.append('a')
```

```
>>> list_2
```

```
[ 'a', 'd', None, 'a' ]
```

```
>>> list_2.remove('a')
```

```
>>> list_2
```

```
[ 'd', None, 'a' ]
```

Python - Collections

Lists

- The `in` operator checks whether a value is contained in a list:

```
>>> 'a' in list_2
```

```
True
```

```
>>> 'c' in list_2
```

```
False
```

Python - Collections

Lists

- The operator `+` is used to concatenate lists:

```
>>> list_1 + list_2
```

```
[ 2, 3, 7, None, 'd', None, 'a' ]
```

- `extend()` adds full lists to a given one:

```
>>> list_1.extend( list_2 )
```

```
>>> list_1
```

```
[ 2, 3, 7, None, 'd', None, 'a' ]
```

Python - Collections

Lists

- A list can be sorted in place using `sort()`:

```
>>> a = [ 7, 2, 5, 1, 3 ]
```

```
>>> a.sort()
```

```
>>> a
```

```
[1, 2, 3, 5, 7]
```

```
>>> b = ['galicia', 'asturias', 'cantabria', 'euskadi', 'navarra']
```

```
>>> b.sort()
```

```
>>> b
```

```
[ 'asturias', 'cantabria', 'euskadi', 'galicia', 'navarra' ]
```

Python - Collections

Lists

- `sort()` accepts a function as sorting key:

```
>>> b = [ 'caladan', 'arrakis', 'corrin', 'ix', 'giedi  
prime' ]
```

```
>>> b.sort( key = len )
```

```
>>> b
```

```
[ 'ix', 'corrin', 'caladan', 'arrakis', 'giedi prime'  
]
```


Python - Collections

Lists

- The `bisect` module manipulates sorted lists using binary search.
- It does not check that a list is actually sorted, using it on unsorted lists yields incorrect results.

```
>>> import bisect
```

```
>>> c = [1, 2, 2, 2, 3, 4, 7]
```

```
>>> bisect.bisect( c, 2 ) # Returns insertion index
```

```
4
```

```
>>> bisect.insort( c, 6 ) # Sorted insertion
```

```
>>> c
```

```
[1, 2, 2, 2, 3, 4, 6, 7]
```

Python - Collections

Slicing

- Sections of indexed collections (such as tuples and lists) can be accessed using slice notation:

```
>>> a = [7, 2, 3, 7, 5, 6, 0, 1]
```

```
>>> a[1:5]
```

```
[2, 3, 7, 5]
```

```
>>> a[3:4] = [6, 3]
```

```
>>> a
```

```
[7, 2, 3, 6, 3, 5, 6, 0, 1]
```

Python - Collections

Slicing

- The slice `[start:stop]` includes the element in the `start` position, but not the element in the `stop` position.
- Any of them may be omitted, in which case the start or ending element of the collection is used by default:

```
>>> a[:5]
```

```
[7, 2, 3, 6, 3]
```

```
>>> a[3:]
```

```
[6, 3, 5, 6, 0, 1]
```

Python - Collections

Slicing

- Negative indices refer to the end of the array:

```
>>> a[-4:]
```

```
[5, 6, 0, 1]
```

```
>>> a[-6:-2]
```

```
[6, 3, 5, 6]
```

Python - Collections

Slicing

- The step can be modified using `[start:stop:step]:`

```
>>> a[::2]
```

```
[7, 3, 3, 6, 1]
```

- A particular use of the step is the reverse a sequence:

```
>>> a[::-1]
```

```
[1, 0, 6, 5, 3, 6, 3, 2, 7]
```

Python - Collections Manipulation

- Oftentimes we want to iterate the elements of a collection and their index at the same time:

```
i = 0
```

```
for x in collection:
```

```
    # do something with x, i
```

```
    i += 1
```

- This is equivalent to

```
for i, x in enumerate(collection):
```

```
    # do something with x, i
```

Python - Collections Manipulation

- `sorted()` returns a sorted list containing the elements of a collection:

```
>>> sorted( [7, 1, 2, 6, 0, 3, 2] )
```

```
[0, 1, 2, 2, 3, 6, 7]
```

```
>>> sorted( 'test string' )
```

```
[' ', 'e', 'g', 'i', 'n', 'r', 's', 's', 't', 't',  
't']
```

```
>>> sorted( set( 'test string' ) )
```

```
[' ', 'e', 'g', 'i', 'n', 'r', 's', 't']
```

Python - Collections

Manipulation

- `zip()` groups the elements in several sequences into a list of tuples:

```
>>> seq1 = ['caladan', 'kaitain', 'giedi prime', 'arrakis']
```

```
>>> seq2 = ['atreides', 'corrino', 'harkonnen']
```

```
>>> zip( seq1, seq2 )
```

```
[('caladan', 'atreides'), ('kaitain', 'corrino'),
```

```
('giedi prime', 'harkonnen')]
```

- The length of the resulting list is given by the length of the shortest input sequence.

Python - Collections Manipulation

- `zip()` is commonly used to simultaneously iterate several sequences:

```
for i, (a,b) in enumerate( zip( seq1, seq2 ) ):
    # do something with i, a, b
```

Python - Collections

Manipulation

- `zip()` can be used to “unzip”:

```
>>> l = [ ('caladan', 'atreides'), ('kaitain', 'corrino'),  
         ('giedi prime', 'harkonnen') ]
```

```
>>> unzip1, unzip2 = zip( *l )
```

```
>>> unzip1
```

```
('caladan', 'kaitain', 'giedi prime')
```

```
>>> unzip2
```

```
('atreides', 'corrino', 'harkonnen')
```

Python - Collections

Manipulation

- `reversed()` builds an iterator over the elements of a sequence in reverse order:

```
>>> reversed( range(10) )
```

```
<listreverseiterator object at 0x7ff664159510>
```

```
>>> list( reversed( range(10) ) )
```

```
[9, 8, 7, 6, 5, 4, 3, 2, 1, 0]
```

Python - Collections

Dictionaries

- A dictionary (`'dict'` type) is an associative array (hashtable).
- It is similar to an array, but indexes an object using a key instead of an integer.

```
>>> empty_dict = {}
```

```
>>> d1 = { 'a': 'a value', 'b': [1,2,3,4] }
```

```
>>> d1
```

```
{ 'a': 'a value', 'b': [1,2,3,4] }
```

Python - Collections

Dictionaries

- The elements of a dictionary can be accessed and inserted using the same syntax as for lists and tuples:

```
>>> d1[7] = 'an integer'
```

```
>>> d1
```

```
{'a': 'a value', 'b': [1,2,3,4], 7: 'an integer'}
```

```
>>> d1['b']
```

```
[1,2,3,4]
```

Python - Collections

Dictionaries

- It can be checked whether a key is contained in a dictionary using `in`:

```
>>> 'b' in d1
```

```
True
```

```
>>> 42 in d1
```

```
False
```

Python - Collections

Dictionaries

- To remove values from a dictionary either `del` or `pop()` can be used:

```
>>> del d1['a']
```

```
>>> d1
```

```
{'b': [1,2,3,4], 7: 'an integer'}
```

```
>>> d1.pop('b')
```

```
[1,2,3,4]
```

```
>>> d1
```

```
{7: 'an integer'}
```

Python - Collections

Dictionaries

- `keys()` and `values()` return the keys and values stored in the dictionary:

```
>>> d1[5] = 'another integer'
```

```
>>> d1.keys()
```

```
[5, 7]
```

```
>>> d1.values()
```

```
['another integer', 'an integer']
```

- The keys and values are not returned in any particular order, but the orders of both lists are consistent (i.e. `d1[d1.keys()[x]] = d1.values()[x]`).

Python - Collections

Dictionaries

- Two dictionaries can be fused using `update()`:

```
>>> d1.update( { 'b': 'caladan', 'c': 'arrakis' } )
```

```
>>> d1
```

```
{ 'c': 'arrakis', 'b': 'caladan', 5: 'another integer',  
7: 'an integer' }
```

Python - Collections

Dictionaries

- A dictionary can be created from two lists:

```
mapping = {}
```

```
for key, value in zip( key_list, value_list ):
    mapping[key] = value
```

- Or:

```
>>> mapping = dict( zip( key_list, value_list ) )
```

Python - Collections

Dictionaries

- In order to be usable as a dictionary key, a Python object must be “hashable”.
- Basic types in Python are hashable, but not mutable containers (e.g., lists).
- A hashable object implements `__hash__()`, `__eq__()`, and `__cmp__()` such that:
 1. The return value of `__hash__()` does not change during the life of the object.
 2. If two objects are equal according to `__eq__()` they must share the same `__hash__()` value.
 3. `__cmp__()` must compare objects consistently.

Python - Collections

Set

- A set is an unsorted collection of unique elements:

```
>>> set( [2,2,2,1,3,3] )
```

```
set([1, 2, 3])
```

```
>>> {2, 2, 2, 1, 3, 3}
```

```
set([1, 2, 3])
```

Python - Collections

Set



Method	Alternate syntax	Description
<code>a.add(x)</code>	--	Add <code>x</code> to set.
<code>a.remove(x)</code>	--	Remove <code>x</code> from set.
<code>a.union(b)</code>	<code>a b</code>	Union of <code>a</code> and <code>b</code> .
<code>a.intersection(b)</code>	<code>a & b</code>	Intersection of <code>a</code> and <code>b</code> .
<code>a.difference(b)</code>	<code>a - b</code>	Set difference.
<code>a.symmetric_difference(b)</code>	<code>a ^ b</code>	Symmetric set difference.
<code>a.issubset(b)</code>	--	True if <code>b</code> is a subset of <code>a</code> .
<code>a.issuperset(b)</code>	--	True if <code>a</code> is a superset of <code>b</code> .
<code>a.isdisjoint(b)</code>	--	True if <code>a</code> and <code>b</code> are disjoint.

Python - Collections

Comprehensions of lists, sets, and dictionaries

- Comprehensions are “syntactic sugar” to generate new collections by operating and filtering preexisting ones.
- The basic syntax is as follows:

```
>>> [expr for val in collection if condition]
```

equals:

```
new_list = []
```

```
for val in collection:
```

```
    if condition:
```

```
        new_list.append( expr )
```

Python - Collections

Comprehensions of lists, sets, and dictionaries

- The filtering condition may be omitted:

```
>>> strings = ['a', 'an', 'the', 'cat', 'car', 'pigeon']
```

```
>>> [x.upper() for x in strings if len(x) > 2]
```

```
['THE', 'CAT', 'CAR', 'PIGEON']
```

```
>>> [x.upper() for x in strings]
```

```
['A', 'AN', 'THE', 'CAT', 'CAR', 'PIGEON']
```

Python - Collections



Comprehensions of lists, sets, and dictionaries

- Using a similar syntax we can write comprehensions of sets and dictionaries:

```
{ key-expr: val-expr for value in collection if condition }
```

```
{ set-expr for value in collection if condition }
```


Python - Collections

Comprehensions of lists, sets, and dictionaries

- We can write nested loops in a comprehension:

```
>>> tuples = ((1, 2, 3), (4, 5, 6), (7, 8, 9))
```

```
>>> [x for tup in tuples for x in tup if x > 3]
```

```
[4, 5, 6, 7, 8, 9]
```

- The order of the loops in a comprehension is the same as in an equivalent code:

```
>>> list = []
```

```
>>> for tup in tuples:
```

```
...     for x in tup:
```

```
...         if x > 3: list.append( x )
```

Python - Collections

Comprehensions of lists, sets, and dictionaries

- It is also valid to nest comprehensions:

```
>>> [[x for x in tup] for tup in tuples]
```

```
[[1, 2, 3], [4, 5, 6], [7, 8, 9]]
```

Python - Functions



- Functions are defined using the reserved word `def`.
- `return` is used to return control to the caller and pass the result.

```
def my_func( x, y, z = 1.5 ):
    if z > 1:
        return z * (x+y)
    else:
        return z / (x+y)
```

Python - Functions



- If the end of the function code is reached without executing any `return` instruction, `None` is returned automatically.
- Functions receive two types of parameters: positional and keyword.
- Keyword arguments are commonly used to provide default values to optional parameters.
- In the example, `x` and `y` are positional parameters, while `z` is a keyword parameter. The function can be called in different ways:

```
>>> my_func( 5, 6,  
z=0.7 )  
6.14
```

```
>>> my_func( 3.14, 7,  
3.5 )  
17.990000000000000002
```

Python - Functions



Name spaces, scopes, and local functions

- A function can access variables in two different scopes: global and local.
- Alternatively, variables can be explicitly defined inside a namespace.
- By default, variables assigned inside a function belong to the local scope.
- The local scope of a function is created when it is called, and initially contains its parameters.
- The local scope is destroyed when the function returns (except for closures, which we will briefly cover later).

Python - Functions

Name spaces, scopes, and local functions



```
def func():  
  
    a = []  
  
    for i in range(5):  
  
        a.append(i)
```

- When `func()` is called, `a` is created. Then, the loop is executed and 5 integers are appended to `a`. Finally, the end of the function body is reached and `a` is destroyed.

Python - Functions

Name spaces, scopes, and local functions



```
>>> a = []
```

```
>>> def func():
```

```
...     for i in range(5):
```

```
...         a.append( i )
```

```
>>> func()
```

```
>>> a
```

```
[0, 1, 2, 3, 4]
```

Python - Functions

Name spaces, scopes, and local functions



```
>>> a = []
```

```
>>> def func():
```

```
...     a = range(5)
```

```
>>> func()
```

```
>>> a
```

```
[]
```


Python - Functions

Name spaces, scopes, and local functions



```
>>> a = []
```

```
>>> def func():
```

```
...     global a
```

```
...     a = range(5)
```

```
>>> func()
```

```
>>> a
```

```
[0, 1, 2, 3, 4]
```

Python - Functions



Name spaces, scopes, and local functions

- New functions can be declared anywhere in the code.
- In particular, it is legal to declare a function nested inside another function. These are called local functions, and are created when the enclosing function is called:

```
def outer_f( x, y, z ):

    def inner_f( a, b, c ):

        pass

    pass
```

- `inner_f()` does not exist until the call to `outer_f()`. As soon as `outer_f()` returns, `inner_f()` is destroyed.
- `inner_f()` can access the variables and functions in the local scope of `outer_f()`, but it cannot add anything to it.

Python - Functions

Returning multiple values

- A notable difference w.r.t. other languages such as C/C++/Java is the ability of a function to return multiple values.

```
def f():
```

```
    a = 5; b = 6; c = 7;
```

```
    return a, b, c
```

```
x, y, z = f()
```

- The implementation is very simple: `f()` is actually returning a single value, which happens to be a tuple.

```
>>> f()
```

```
(5, 6, 7)
```

Python - Functions

Functions as objects

- A Python function is just a special type of object which defines the `()` operator.
- As such, it is possible to use functors (pointers to function objects) to code complex operations in a simple way.
- E.g., we want to perform cleanup operations on the strings in the following array:

```
>>> planets = [ '    Caladan ', 'Ix!', 'Ix', 'ix',  
'aRraKIs', 'giedi prime##', 'Salusa secundus?' ]
```

- We want to build a list of uniform strings for its analysis. We need to apply removal of unnecessary spaces and symbols, and to fix capitalization.

Python - Functions

Functions as objects



```
import re # Regular Expression module

def clean_strings( strings ):

    result = []

    for value in strings:

        value = value.strip()

        value = re.sub('[!#?]', '', value) # removes !#? symbols

        value = value.title()

        result.append( value )

    return result
```

Python - Functions

Functions as objects

```
def clean_strings_v2( strings, ops ):  
  
    result = []  
  
    for value in strings:  
  
        for function in ops:  
  
            value = function( value )  
  
        result.append( value )  
  
    return result  
  
  
def remove_punctuation( value ):  
  
    return re.sub( "[!#?]", "", value )  
  
  
clean_ops = [ str.strip, remove_punctuation, str.title ]
```

Python - Functions

Functions as objects

- It is possible to pass functions as arguments to other functions.
- E.g., Python provides mechanisms to apply a function to a list of objects:

```
>>> map( str.strip, planets )
```

```
['Caladan', 'Ix!', 'Ix', 'ix', 'aRraKIs', 'giedi prime##',  
'Salusa secundus?']
```

Python - Functions

Anonymous (λ) functions

- An anonymous function, or lambda function, is a single-instruction functional expression the result of which is its return value:

```
def short_function(x):  
    return x*2
```

```
equiv_anon = lambda x: x*2
```


Python - Functions

Anonymous (λ) functions

- Oftentimes it is simpler to use a reference to a lambda function than to write an ad-hoc, named one:

```
>>> a = [4, 0, 1, 5, 6]
```

```
>>> map( lambda x: x*2, a )
```

```
[8, 0, 2, 10, 12]
```

```
>>> strings = ['caladan', 'ix', 'corrin', 'giedi prime']
```

```
>>> strings.sort( key = lambda x: len(set(list(x))) )
```

```
# Sort by number of different letters in the string
```

```
>>> strings
```

```
['ix', 'caladan', 'corrin', 'giedi prime']
```

Python - Functions

Closures

- A closure is a dynamically-generated function which is returned by another function.
- The distinguishing characteristic of a closure is that it is capable of accessing the local scope of its generating function after the latter returns.

Python - Functions

Closures

```
>>> def make_closure( a ):
...     def closure():
...         print( "Variable in the local scope: %d" % a )
...     return closure

>>> closure = make_closure(5)

>>> closure()

'Variable in the local scope: 5'
```

Python - Functions

Closures

- In the example, a closure was created with an immutable internal state (the integer `a`).
- Mutable variables can also be used in a closure. These can be modified, dynamically altering the behavior of the closure.

Python - Functions

Closures

```
>>> def make_watcher():  
  
...     have_seen = set([])  
  
...     def has_been_seen( x ):  
  
...         if x in have_seen: return True  
  
...         else: have_seen.add(x)  
  
...         return False  
  
...     return has_been_seen  
  
>>> watcher = make_watcher()  
  
>>> vals = [5, 6, 1, 5, 1, 6, 3, 5]  
  
>>> [watcher(x) for x in vals]  
  
[False, False, False, True, True, True, False, True]
```

Python - Functions

Closures

- The local variables of a closure can be **modified**.
- No new variables can be **added** to the scope of a closure. A workaround is to add key/value pairs to a dictionary in the scope.
- Closures allow to build generic functions with plenty of options, that can be dynamically instantiated into specialized, efficient and simple variants.

Python - Functions

Closures



```
>>> def format_and_pad( template, space ):  
...     def formatter( x ):  
...         return (template % x).rjust( space )  
...     return formatter  
  
>>> fmt = format_and_pad( "%.4f", 15 )  
  
>>> fmt(1.756)  
  
'          1.7560'
```

Python - Functions

Extended syntax: `*args`, `**kwargs`

- Functions are called using a mix of positional and keyword parameters:

```
>>> func( a, b, c, d = d_value, e = e_value )
```

- Internally, this function:
 1. Receives an `args` tuple containing its positional parameters.
 2. Receives a `kwargs` dictionary containing its keyword parameters:
 3. Performs the following assignment:

```
>>> (a, b, c) = args
```

```
>>> d = kwargs.get( 'd', d_default_value )
```

```
>>> e = kwargs.get( 'e', e_default_value )
```




Python - Functions

Extended syntax: *args, **kwargs

```
def g( x, y, z=1 ): return (x+y) / z

def hello_world_then_call( f, *args, *kwargs ):

    print 'args is', args

    print 'kwargs is', kwargs

    print "Hello world! Now I'm going to call %s" % f

    return f( *args, *kwargs )

>>> hello_world_then_call( g, 1, 2, z=5 )

args is (1, 2)

kwargs is {'z': 5.0}

Hello world! Now I'm going to call <function g at 0x2dd5cf8>

0.6
```

Python - Functions



Partial function application

- Partial function application consists in creating new functions from preexisting ones by fixing some of their parameters:

```
>>> def add(x, y): return x+y
```

```
>>> add_5 = lambda y: add(5,y)
```

- Alternatively:

```
>>> from functools import partial
```

```
>>> add_5 = partial( add, 5 )
```

Python - Functions

Generators

- A generator is a function which returns a sequence of values in a lazy way, stopping its execution after each value in the sequence.
- Generators are useful to generate large, iterable sequences in a memory-efficient way (e.g., `range()` vs `xrange()`).
- Generators are declared as a function which returns a value using `yield` instead of `return`.

Python - Functions

Generators

```
>>> def squares( n = 10 ):
...     for i in xrange( 1, n+1 ):
...         print "Generating squares from 1 to
%d"%(n**2)
...         yield i ** 2
>>> gen = squares()
>>> gen
<generator object squares at 0x7fd9e3e796e0>
```

Python - Functions

Generators

- When the generator function is called, no code is executed.
- Each element must be explicitly requested:

```
>>> for x in gen:
```

```
...     print x
```

```
Generating squares from 1 to 100
```

```
1
```

```
Generating squares from 1 to 100
```

```
4
```

```
.
```

```
.
```

```
.
```

Python - Functions

Generators

```
def make_change( amount, coins=[1, 2, 5, 10, 20, 50], hand=[] ):

    if amount == 0: yield hand

    for coin in coins:

        if coin > amount or (len(hand) > 0 and hand[-1] < coin):
            continue

        for result in make_change(amount-coin, coins=coins,
            hand=hand+[coin]):

            yield result
```

```
>>> len(list(make_change(53)))
```

Python - Functions

Generator expressions

- A compact way to declare a simple generator is to use a generator expression:

```
>>> gen = ( x ** 2 for x in xrange(100) )
```

```
>>> gen
```

```
<generator object <genexpr> at 0x7fd9e3e798c0>
```

- This expression is equivalent to:

```
def gen():  
    for x in xrange(100):  
        yield x**2
```

Python

Files and Operating System

- To open a file, `open()` is called passing a relative or absolute path:

```
>>> f = open( 'folder/file.txt' )
```

- By default, the file is opened in read-only mode (`'r'`). The file can be seen as a collection of lines, and iterated using a `for` loop:

```
>>> for line in f:
```

```
...     # do something with line
```


Python

Files and Operating System



Opening mode	Description
<code>r</code>	Read only.
<code>w</code>	Write only. Creates a new file, overwriting any previous one.
<code>a</code>	Concatenate to a file (created if it does not exist).
<code>r+</code>	Read-write.
<code>b</code>	Binary mode (usage example: <code>'rb'</code>).
<code>U</code>	Use universal end-of-line mode. Translates any end-of-line marker in the file to <code>'\n'</code> .

Python

Files and Operating System

- To write to a file we can use the methods `write()` or `writelines()`:
- E.g., to remove white lines from a file:

```
>>> f_in = open( path, 'r' )
```

```
>>> f_out = open( 'tmp.txt', 'w' )
```

```
>>> f_out.writelines( [x for x in f_in if len(x) > 1]
)
```

Python

Files and Operating System



Method	Description
<code>read([size])</code>	Reads data from the file as a string. The optional argument <code>[size]</code> is the number of bytes to read. Without it, the entire file is read.
<code>readlines([size])</code>	Same as <code>read()</code> , but returns a list of strings (one per line in the file). Without <code>[size]</code> the entire file is read.
<code>write(str)</code>	Writes <code>str</code> to the file.
<code>writelines(str)</code>	Writes a list of strings to the file, one per line.

Python

Files and Operating System



Method	Description
<code>close()</code>	Closes the file.
<code>flush()</code>	Synchronizes the I/O buffer to disk.
<code>seek(pos)</code>	Moves the file pointer to <code>pos</code> .
<code>tell()</code>	Returns the current position of the file pointer.
<code>closed</code>	<code>True</code> if the file is closed.

NumPy



- NumPy (***Numerical Python***) is a fundamental package that enables efficient array and vector processing.
- It provides:
 - An `ndarray` class, which represents a multidimensional array and provides *vectorized* arithmetic operations and *broadcasting* capabilities.
 - *Vectorized* operations are applied to whole arrays (without using Python loops).
 - Operations for lineal algebra, randomness, signal processing, ...
 - Tools for the effective integration of Python and C/C++/Fortran.

NumPy - ndarray

002-Numpy.ipynb

- n-dimensional array object.
- Allows to perform operations on large data blocks with scalar syntax.
- Unlike lists, a NumPy array is usually homogeneous: it holds objects of a specific, predefined type.
- An `ndarray` contains, among others, the following two properties:
 - `shape`: tuple containing the array dimensions.
 - `dtype`: instance of the `dtype` class which specifies element datatype.

NumPy - ndarray Creation

002-Numpy.ipynb

Function	Description
<code>array()</code>	Builds a new array from the input sequence (e.g., a list). Copies input data to the new array.
<code>asarray()</code>	Converts the input sequence to an array, but it does not copy data if the input is already an <code>ndarray</code> .
<code>arange()</code>	Same as <code>range()</code> but returning an <code>ndarray</code> instead of a list.
<code>ones()</code> <code>ones_like()</code>	Build an array containing all 1s from a tuple specifying the desired dimensions, or copying <code>shape</code> and <code>dtype</code> from another array.

NumPy - ndarray Creation

002-Numpy.ipynb

Function	Description
<code>zeros()</code> <code>zeros_like()</code>	Build an array containing all 0s from a tuple specifying the desired dimensions, or copying <code>shape</code> and <code>dtype</code> from another array.
<code>empty()</code> <code>empty_like()</code>	Build an array without initializing its data from a tuple specifying the desired dimensions, or copying <code>shape</code> and <code>dtype</code> from another array.
<code>eye()</code> , <code>identity()</code>	Build an array with 1s in its main diagonal and 0s elsewhere. <code>identity()</code> requires square dimensions.

NumPy - ndarray Datatypes



Type	Code	Description
<code>int8,</code> <code>uint8</code>	<code>i8,</code> <code>u8</code>	Signed/unsigned integer, 8 bits (1 byte).
<code>int16,</code> <code>uint16</code>	<code>i16,</code> <code>u16</code>	Signed/unsigned integer, 16 bits (2 bytes).
<code>int32,</code> <code>uint32</code>	<code>i32,</code> <code>u32</code>	Signed/unsigned integer, 32 bits (4 bytes).
<code>int64,</code> <code>uint64</code>	<code>i64,</code> <code>u64</code>	Signed/unsigned integer, 64 bits (8 bytes).

NumPy - ndarray

Datatypes



Type	Code	Description
<code>float16</code>	<code>f2</code>	Floating point, half precision (2 bytes).
<code>float32</code>	<code>f4</code> <code>o</code> <code>f</code>	Floating point, single precision (4 bytes). Compatible with a C <code>float</code> .
<code>float64</code>	<code>f64</code> <code>o</code> <code>d</code>	Floating point, double precision (8 bytes). Compatible with a C <code>double</code> and with a Python <code>float</code> .
<code>float128</code>	<code>f16</code> <code>o</code> <code>g</code>	Floating point, extended precision (16 bytes).

NumPy - ndarray Datatypes



Type	Code	Description
<code>complex64</code>	<code>c8</code>	Complex number represented as 2 x float32.
<code>complex128</code>	<code>c16</code>	Complex number represented as 2 x float64.
<code>complex256</code>	<code>c32</code>	Complex number represented as 2 x float128.

NumPy - ndarray

Datatypes



Type	Code	Description
<code>bool</code>	<code>?</code>	Boolean type (<code>True</code> or <code>False</code>).
<code>object</code>	<code>O</code>	Python object (reference).
<code>string_</code>	<code>S_</code>	String type with fixed length (1 byte per character). E.g., <code>S10</code> represents a 10-character string.
<code>unicode_</code>	<code>U_</code>	Unicode type with fixed length (the number of bytes per character varies per platform). Similar to <code>string_</code> .

NumPy - ndarray Datatypes

002-NumPy.ipynb

- `ndarray` can be cast to a different datatype using `astype()`.
- The call to `astype()` always creates a new copy of the array memory, even if the data is cast to its current type.
- If the cast fails for any reason, a `TypeError` exception is raised.

NumPy - ndarray

Broadcasting

- When two arrays are operated, NumPy must decide what operation to perform in the first place.
- If both arrays have the same dimensionality, the operation is performed in an element-wise fashion.
- Otherwise, NumPy tries to *broadcast* (replicate) the “smallest” array to match it to the largest one.

NumPy - ndarray

Broadcasting

- The broadcasting process begins by the innermost dimension (the rightmost ones), and moves towards the outer ones (to the left).
- Two dimensions are compatible if:
 1. Both have the same number of elements, or
 2. One of them has a single element.
- If none of the previous conditions is met, the operation raises a `ValueError`.
- If the second condition is met, the array with a single element in that particular dimension is copied `n` times to match both.

NumPy - ndarray

002-NumPy.ipynb

Broadcasting

- The arrays are not required to have the same number of dimensions. Non-existing dimensions are assumed to have a single element.
- For example:

A	:	256	x	256	x	3	A	:	256	x	256	x	3
B	:					3	B	:	256	x	1	x	3
A*B:		256	x	256	x	3	A*B:		256	x	256	x	3

NumPy - ndarray

002-NumPy.ipynb

Basic indexing

- 1-dimensional arrays behave in a similar way to lists.
- An important difference is that in NumPy a slice is not a copy, but a *view* over the array data. This improves performance, but also means that slice operations may have collateral effects on the original array.
- It is possible to obtain a copy of a slice (or any array) by calling the `copy()` method of `ndarray`.

NumPy - ndarray

Basic indexing

002-NumPy.ipynb

The elements in each index of an n-dimensional array are not scalars, but (n-1)-dimensional arrays.

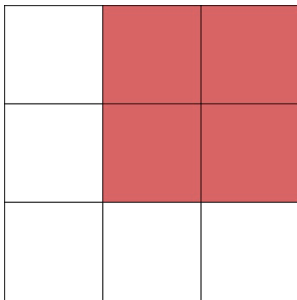
		axis 1		
		0	1	2
axis 0	0	0, 0	0, 1	0, 2
	1	1, 0	1, 1	1, 2
	2	2, 0	2, 1	2, 2

NumPy - ndarray

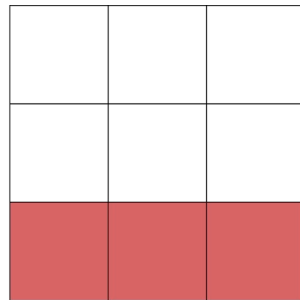
002-NumPy.ipynb

Indexing with slices

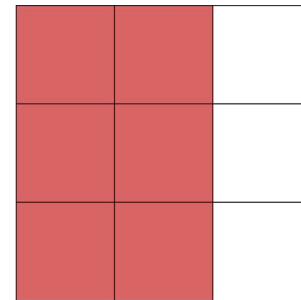
- The slice syntax is also valid for indexing arrays.
- Working with n-dimensional arrays, it is possible to use any combination of basic indexing and slicing.



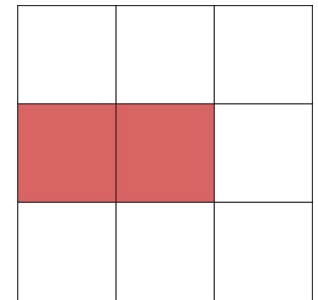
`arr[:2, 1:]`
shape = (2,2)



`arr[2]`
shape = (3,)



`arr[:, :2]`
shape = (3,2)



`arr[1:2, :2]`
shape = (1,2)

NumPy - ndarray

002-NumPy.ipynb

Boolean indexing

- It is possible to index an array using a boolean array.
- The boolean array must have the same dimensionality as the axis it must index.
- Boolean indexing can be used to select elements in an array using data semantics (e.g., the elements which are greater than a particular number). The conditions do not need to refer to the array being indexed.
- Boolean indexing always returns a copy of the data (not a view).
- When accessing an n-dimensional array, boolean indexing can be mixed together with basic and sliced indexing.

NumPy - ndarray

002-NumPy.ipynb

Fancy indexing

- The term *fancy indexing* describes the indexing of an array using an integer collection containing the indices of the elements to select.
- If a list of lists is provided, the behavior changes: a 1D array is returned, containing the data indexed by the tuples resulting from applying `zip()` to the input lists.
- Fancy indexing always returns a copy of the data, not a view.
- When accessing an n-dimensional array, fancy indexing can be mixed together with basic, sliced, and boolean indexing.

NumPy - ndarray

002-NumPy.ipynb

Transposition

- `ndarray` provides the `transpose()` method and the `T` attribute.
- For n-dimensional arrays, `transpose()` accepts as input a tuple specifying an arbitrary dimensional permutation.
- The `T` attribute is a quick way of accessing “classic” transposition, in which a matrix is flipped over its main diagonal.
- The `swapaxes()` is used to swap any two axes.
- All these operations return a view over the original array.

NumPy - Universal functions

002-NumPy.ipynb

- A universal function, or `ufunc`, performs element-wise operations over an `ndarray`.
- The term describes a wrapper over a simple function which reads one or more scalars and returns a single scalar.
- Many `ufuncs` take a single element as input (e.g., `sqrt()` or `exp()`). These are *unary* `ufuncs`.
- Other `ufuncs` take two elements and return one (*binary* `ufuncs`, e.g., `add()` or `maximum()`).
- Some `ufuncs` return more than one array (e.g., `modf()` returns the integral and fractional parts of a floating point array).

NumPy - Universal functions



Unary ufuncs

Function	Operation
<code>abs, fabs</code>	Absolute value of integers, floating point or complex numbers. <code>fabs()</code> is a faster version for non-complex numbers.
<code>sqrt</code>	Square root: <code>arr ** 0.5</code> .
<code>square</code>	Square: <code>arr ** 2</code> .
<code>exp</code>	Exponentiation: <code>np.e ** arr</code> .
<code>log, log10, log2, log1p</code>	Natural, base 10, and base 2 logarithms. <code>log1p(x)</code> returns <code>log(1+x)</code> .

NumPy - Universal functions



Unary ufuncs

Function

Return value

<code>sign</code>	The sign of each element in the array: 1 (positive), 0 (zero), -1 (negative).
<code>ceil / floor / rint</code>	Round up / down / to closer integer.
<code>modf</code>	Integral and fractional parts of an array.
<code>isnan</code>	True if an element is <code>np.nan</code> (NaN).
<code>isfinite, isinf</code>	Boolean array indicating if an element is finite / infinite.

NumPy - Universal functions



Unary ufuncs

Function	Return value
<code>cos, cosh, sin, sinh, tan, tanh</code>	Trigonometric and hyperbolic functions.
<code>arccos, arccosh, arcsin, arcsinh, arctan, arctanh</code>	Inverse trigonometric and hyperbolic functions.
<code>logical_not</code>	Element-wise boolean NOT.

NumPy - Universal functions



Binary ufuncs

Function	Return value
<code>add</code>	Addition of two input arrays.
<code>subtract</code>	Subtracts the second parameter from the first.
<code>multiply</code>	Multiplication of two input arrays.
<code>divide / floor_divide</code>	Division / integral division of the input.
<code>power</code>	Exponentiation of the bases in the first array to the exponents in the second array.

NumPy - Universal functions



Binary ufuncs

Function	Return value
<code>maximum / fmax</code>	Element-wise maximum. <code>fmax()</code> ignores NaN values.
<code>minimum / fmin</code>	Element-wise minimum. <code>fmin()</code> ignores NaN values.
<code>mod</code>	Element-wise modulo.
<code>copysign</code>	Copies the signs of the element in the second array to the values in the first.

NumPy - Universal functions



Binary ufuncs

Function	Return value
<code>greater / greater_equal / less / less_equal / equal / not_equal</code>	Boolean operations comparing whether each element in the first array is greater / greater or equal / less / less or equal / equal / different than each element in the second array.
<code>logical_and / logical_or / logical_xor</code>	Element-wise boolean operations AND / OR / XOR.

NumPy - Data processing

002-NumPy.ipynb

- Using NumPy we can express data processing algorithms that would otherwise require loops using array operations.
- This *vectorization* process usually comes together with significant performance improvements.
- We will explore different ways of using NumPy to process data.

NumPy - Data processing

Conditional logic

002-NumPy.ipynb

- `numpy.where()` allows to write a vectorized version of the expression `x if c else y`.
- Using pure Python, this can be written as:

```
>>> result = [(x if c else y) for x, y, c in zip(xarr, yarr, cond)]
```

- However, this version...
 - Is very slow.
 - works with 1-dimensional arrays only.
- The equivalent `where()` version is written:

```
>>> result = np.where( cond, xarr, yarr )
```

NumPy - Data processing

Mathematics and statistics



Method	Description
<code>sum</code>	Summation of all the elements of an array along a given optional axis.
<code>mean</code>	Arithmetic mean.
<code>std, var</code>	Standard deviation and variance.
<code>min, max</code>	Minimum and maximum.

NumPy - Data processing

Mathematics and statistics



002-NumPy.ipynb

Method	Description
<code>argmin, argmax</code>	Indices of minimum and maximum elements.
<code>cumsum</code>	Cumulative sum.
<code>cumprod</code>	Cumulative product.

NumPy - Data processing

Boolean arrays

002-NumPy.ipynb

- When working with boolean arrays, the previous methods use `True == 1` and `False == 0`. This allows to reduce boolean values using `cumsum()` / `cumprod()` instead of **AND** / **OR**.
- There are however ad-hoc methods for reducing boolean arrays: `any()` and `all()`.
- These methods work with non-boolean arrays by interpreting `non-zero == True`, `zero == False`.

NumPy - Data processing

Sorting

002-NumPy.ipynb

- As with lists, NumPy arrays can be sorted in place using `sort()`.
- `sort()` accepts an optional `axis` parameter which indicates the axis to sort over in multidimensional arrays.
- `numpy.sort()` returns a sorted copy of the array.

NumPy - Data processing

Set logic



Function	Description
<code>unique(x)</code>	Computes the unique elements in array <code>x</code> (1D).
<code>intersect1d(x, y)</code>	Common elements in <code>x</code> and <code>y</code> .
<code>union1d(x, y)</code>	Union of <code>x</code> and <code>y</code> .
<code>in1d(x, y)</code>	Boolean array indicating whether each element of <code>x</code> is in <code>y</code> .
<code>setdiff1d(x, y)</code>	Set difference, <code>x-y</code> .
<code>setxor1d(x, y)</code>	Symmetric set difference.

NumPy - Array I/O

Storing arrays to disk

- `np.save()` and `np.load()` store arrays to disk in binary format.
- By default, the data is stored uncompressed with extension `.npy`.

```
>>> np.save( path_string, arr )
```

```
.
```

```
.
```

```
.
```

```
>>> arr = np.load( path_string )
```

NumPy - Array I/O

Storing arrays to disk

- We can store several arrays to the same file compressed using ZIP and with `.npz` extension using `np.savez()`:

```
>>> np.savez( path_string, a = arr1, b = arr2 )
```

- Reading a `.npz` file NumPy returns a dictionary object that loads individual arrays in a lazy way:

```
>>> arch = np.load( path_string )
```

```
>>> arch['b']
```

```
array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
```

NumPy - Array I/O

Storing arrays to disk

- It is often useful to read arrays stored in text format.
- NumPy provides `np.loadtxt()` and `np.genfromtxt()`.
- Both can be adapted to a specific input format, varying comments format, delimiters, conversion functions, rows to be skipped, or columns to parse.
- `np.savetxt()` performs the reverse operation: it writes an array to a CSV-like file.
- `np.genfromtxt()` is similar to `np.loadtxt()`, but it supports treatment of missing values and structuring of the output array.

NumPy - Linear algebra

(`numpy.linalg`)



Function

Description

<code>diag</code>	Returns the elements in the diagonal of a square matrix as a 1D array, or builds a square matrix from a 1D array.
<code>dot</code>	Matrix product.
<code>trace</code>	Sum of the elements in the matrix diagonal.
<code>det</code>	Determinant.
<code>eig</code>	Eigenvalues and eigenvectors.

NumPy - Linear algebra

(`numpy.linalg`)



Function

Description

<code>inv</code>	Inverse of an square matrix.
<code>pinv</code>	Moore-Penrose pseudo-inverse.
<code>qr</code>	QR decomposition.
<code>svd</code>	Singular-value decomposition.
<code>solve</code>	Solves $Ax = b$, with A a square matrix.
<code>lstsq</code>	Least squares solution to $y = Xb$.

NumPy - Random numbers

(`numpy.random`)



Function	Description
<code>seed</code>	Changes the generator seed.
<code>permutation</code>	Random permutation of a sequence.
<code>shuffle</code>	Random permutation of a sequence (in-place).
<code>rand</code>	Samples a uniform distribution.
<code>randint</code>	Samples integers from a uniform distribution.

NumPy - Random numbers

(`numpy.random`)



Function	Description
<code>randn</code>	Samples a standard normal distribution.
<code>binomial</code>	Samples a binomial distribution.
<code>normal</code>	Samples a normal distribution.
<code>beta</code>	Samples a beta distribution.

NumPy - Random numbers

(`numpy.random`)

002-NumPy.ipynb

Function	Description
<code>chisquare</code>	Samples a χ^2 distribution.
<code>gamma</code>	Samples a gamma distribution.
<code>uniform</code>	Samples a uniform $[0, 1)$ distribution.

Matplotlib



- The previous examples included some plotting, but without going into details.
- Matplotlib provides mostly 2D plotting capabilities, although it includes limited 3D functionalities.
- The basic supported plot types are lines, bar charts, histograms, pie charts, and variations of them.
- This section briefly introduces Matplotlib. It has a vast amount of options that cannot be covered in detail. We will later focus on higher level libraries instead.

Matplotlib - Introduction



- The `%matplotlib` magic command automatically configures iPython to show plots.
- By default, iPython detects the proper backend for the current window manager. Using `%matplotlib inline` the plots will be shown embedded inside a QT or Notebook environment.
- A basic Matplotlib plot includes the following elements:
 - `x` and `y` axes: horizontal and vertical axes, respectively.
 - Tick marks in the `x` and `y` axes.
 - Tick labels, showing axis values.
 - Drawing area, called canvas.

Matplotlib - Introduction

003-Matplotlib.ipynb

- `plot()` is used to draw lines and marks.
- It accepts pairs of `x` and `y` sequences, which must have the same lengths, together with a string indicating how to plot the data.
- Adjacent points are joined using straight lines.
- Points and lines are drawn following the requested style.
- The function returns a list of the lines which have been added to the current figure.
- If only one sequence is passed to the function, it is assumed to contain the values of the `y` axis. Values for the horizontal axis will be automatically generated as `x = range(len(y))`.

Matplotlib - Introduction

003-Matplotlib.ipynb

Function	Description
<code>figure</code>	Creates a new figure. Accepts an integer that acts as a unique identifier for this figure. This integer can be used to programatically change the active figure.
<code>subplot(x, y, z)</code>	Divides a figure into a mesh of subfigures with <code>x</code> rows and <code>y</code> columns. Besides, it activates subfigure number <code>z</code> ($1 < z < x * y$, row major order).

Matplotlib - Introduction

003-Matplotlib.ipynb

Function

Description

<code>bar / barh</code>	Create vertical / horizontal bar charts. To plot stacked bars the <code>bottom</code> parameters is used, indicating the starting point for the new bars.
<code>boxplot</code>	Box plots.
<code>scatter</code>	Draws points, but it does not join them with lines unlike <code>plot()</code> .
<code>hist</code>	Histograms.
<code>pie</code>	Pie charts.

Matplotlib - Design

003-Matplotlib.ipynb

Function

Description

<code>title</code>	Plot title. Like all other text-related Matplotlib functions, it accepts LaTeX syntax.
<code>xlim / ylim</code>	Configures the limits of the x / y axes.
<code>autoscale</code>	Automatic axes limits.
<code>xticks / yticks</code>	Configures the ticks for each axis and, optionally, the tick labels.

Matplotlib - Design

003-Matplotlib.ipynb

Function	Description
<code>axis</code>	Configures both axes' limits.
<code>axhline</code> / <code>axvline</code>	Draws a horizontal / vertical line,
<code>axhspan</code> / <code>axvspan</code>	Draws rectangles that cover the entire width / height of the plot.

Matplotlib - Design

Axes

003-Matplotlib.ipynb

- A reference to an instance of `matplotlib.axes.Axes` allows to configure plots with a high level of detail.
- Such reference can be obtained by calling `matplotlib.pyplot.gca()` (Get Current Axes).
- E.g., calling `ax.xaxis.set_major_locator(matplotlib.ticker.MultipleLocator(10))` sets ticks in the x axis of the `ax` axes in values multiple of 10.
- E.g., calling `ax.xaxis.set_major_formatter()` allows to specify functions that will be used to format the tick labels in the x axis of the `ax` axes.

Matplotlib - Design

Legends and annotations

003-Matplotlib.ipynb

- All plotting functions (`plot()`, `hist()`, etc.) accept a `label` parameter indicating the name to use in the legend.
- `matplotlib.pyplot.legend()` automatically builds a legend from a list of handles. If no such list is provided, all the elements in the figure with a label not starting with “_” are listed.
- Otherwise, we can provide the lists of lines, etc. returned by each plotting function to specify which elements in the figure should be added to the legend.
- `matplotlib.pyplot.annotate()` adds a textual annotation and the specified place in the figure.

Matplotlib - Styles

003-Matplotlib.ipynb

Property	Value type	Description
<code>alpha</code>	<code>float</code>	Transparency used for the element.
<code>color</code>	Matplotlib color	Line / marker color.
<code>dashes</code>	Sequence	Line pattern.
<code>label</code>	<code>string</code>	Text to label this element in the plot legend.
<code>linestyle</code>	See docs	Line type.
<code>linewidth</code>	<code>float</code>	Line width, in points.

Matplotlib - Styles

003-Matplotlib.ipynb

Property	Value type	Description
marker	See docs	Marker used for points in a line.
mec	Matplotlib color	Marker Edge Color.
mew	float	Marker Edge Width.
mfc	Matplotlib color	Marker Face Color.
markersize	float	Marker size, in points.

Matplotlib - Styles

003-Matplotlib.ipynb

Property	Value type	Description
<code>solid_capstyle</code>	<code>['butt' 'round' 'projecting']</code>	End of line style.
<code>solid_joinstyle</code>	<code>['miter' 'round' 'bevel']</code>	Line union style.
<code>visible</code>	<code>[True False]</code>	Visibility.
<code>xdata</code>	<code>numpy.ndarray</code>	X axis data.
<code>ydata</code>	<code>numpy.ndarray</code>	Y axis data.
<code>Zorder</code>	Número	Stacking order in the Z axis.

Matplotlib - Colors

003-Matplotlib.ipynb

- The `matplotlib.colors` module includes utilities for defining and converting colors.
- Basic predefined colors can be referenced using a single letter: b (*blue*), g (*green*), r (*red*), c (*cyan*), m (*magenta*), y (*yellow*), k (*black*), w (*white*).
- Shades of gray can be codified using a floating point number in `[0, 1]`.
- Other colors can be specified using different formats:
 - HTML hex string: `"#eeeeff"`
 - (R, G, B) tuple, with R, G and B in `[0, 1]`.
 - HTML string: `"red"`, `"burlywood"`, `"chartreuse"`, ...

Matplotlib - Colors

Colormaps

003-Matplotlib.ipynb

- The `matplotlib.cm` module includes a set of colormaps to use with plots.
- Colormaps are useful when using `imshow()`, which is similar to `plot()` but interprets the values in the input array as indices to a colormap. If the input array is 2-dimensional, an image will be plotted.

Matplotlib - Saving to file

003-Matplotlib.ipynb

- Figures can be saved to a file using `matplotlib.pyplot.savefig()`.
- The file format to use is inferred from the extension in the provided path.
- The most relevant parameters that control the quality of the stored figure are:
 - `dpi`: dots per inch.
 - `bbox_inches`: inches of whitespace surrounding the figure.
- The figure does not need to be stored to a file: it can be written to any object which supports I/O, such as `StringIO`.

Matplotlib - 2D histograms



Parameter to <code>hist()</code>	Description
<code>bins</code>	Number of classes to use in the histogram, or sequence containing the boundaries between classes.
<code>range</code>	Range for each class. If <code>bins</code> is not used, this can be provided as a sequence.
<code>normed</code>	If <code>True</code> , values are normalized and the result is a probability density function.

Matplotlib - 2D histograms

003-Matplotlib.ipynb

Parameter to `hist()`

Description

<code>histtype</code>	<p>By default, bar chart. Other values:</p> <ul style="list-style-type: none">• <code>barstacked</code>: stacks bars when working with multiple data series.• <code>step</code>: line without filling.• <code>stepfilled</code>: line with filling.
<code>align</code>	How to align bars for each class: <code>mid</code> , <code>left</code> , <code>right</code> .
<code>color</code>	Colors to be used.
<code>orientation</code>	<code>horizontal</code> or <code>vertical</code> .

Matplotlib - Pie charts

003-Matplotlib.ipynb

Parameter to `pie()`

Description

<code>explode</code>	Fraction of the pie ratio to use as offset for each slice.
<code>autopct</code>	String or function indicating how to label each slice with a numeric value.
<code>pctdistance</code>	Ratio between the center of each slice and the start of the <code>autopct</code> text.
<code>labeldistance</code>	Radial distance to draw labels.
<code>startangle</code>	Rotation angle of the origin.
<code>wedgeprops</code>	Dictionary containing slice properties.

Matplotlib - 3D plots

003-Matplotlib.ipynb

- Although Matplotlib is focused in 2D plots, there are several toolkits provided limited 3D capabilities.
- The `mpl_toolkits.mplot3d` module provides methods to create 3D scatter plots, surfaces, lines and meshes. Its interface is very similar to the original Matplotlib.
- Main difference: axes are instances of `mpl_toolkits.mplot3d.Axes3D`. Projections are performed by specialized 3D classes, while other parts of the figure (labels, ticks, etc.) are managed by vanilla Matplotlib.

Bibliography



- *Learning Python*. Mark Lutz. O'Reilly Media, 2013 (5ª edición).
- *Programming Python*. Mark Lutz. O'Reilly Media, 2011 (4ª edición).

Data Analytics: Pandas

- Intro to Pandas.
- Data I/O.
- Data wrangling.
- Visualization.
- Aggregation.
- Time series.

Pandas



- Pandas provides data structures and methods to improve the structured data processing capabilities of native Python.
- The basic data structure in Pandas is the `DataFrame` (similar to `data.frame` in R). It is a 2D table, conceptually similar to an Excel spreadsheet.
- Pandas combines the array processing capabilities of NumPy with the flexibility of spreadsheets and relational databases.
- Provides indexing, reshaping, splitting, aggregation, and selection of subsets of data.

Pandas



- Pandas design objective is to provide new data management capabilities to the Python ecosystem:
 - Data structures with indexed axes supporting explicit or implicit data alignment.
 - Seamless processing of time series (timestamp-indexed data).
 - Arithmetic operations and reductions along axes.
 - Flexible management of unknown (null) data.
 - Merges and similar operations typical of relational databases.

Pandas - Data structures

005-Pandas-Intro.ipynb

Series

- A `Series` object represents an object conceptually similar to a 1D array.
- It actually contains two different arrays: a data array and an index array, which labels data.
- Both arrays are essentially NumPy arrays with their own datatypes.
- Series are similar to dictionaries, as they can be seen as key-value pairs. In fact, a constructor is provided to build a `Series` object from a `dict`.

Pandas - Data structures

005-Pandas-Intro.ipynb

DataFrame

- A `DataFrame` object represents a tabular structure, similar to a spreadsheet (or a `data.frame` in R).
- It contains an ordered collection of columns, each of which may have a different datatype (numeric, string, etc.).
- Includes an index on its rows and another one on its columns. In this sense, it can be seen as an aggregation of series, all of them sharing the same index.
- Internally, the data is stored in a 2D format (bidimensional NumPy array), although higher dimensional data may be represented using hierarchical indices.

Pandas - Data structures



DataFrame () constructor

Parameter (type)	Description
<code>2D ndarray</code>	Data array, with optional row and column labels.
<code>dict of arrays, lists, or tuples</code>	Each sequence becomes a column of the <code>DataFrame</code> . All of them might have the same length.
<code>dict of Series</code>	Each value becomes a column. The indices in each series are unified if an explicit index is not provided.

Pandas - Data structures



DataFrame () constructor

Parameter (type)	Description
<code>dict of dicts</code>	Each internal dictionary becomes a column. The keys in the different dictionaries are unified as in a <code>dict of Series</code> .
<code>list of dicts or Series</code>	Each item becomes a row in the <code>DataFrame</code> . The union of the keys of the dictionaries or the indices of the series is used for the column labels.
<code>list of lists or tuples</code>	Same as providing an 2D <code>ndarray</code> .

Pandas - Data structures



DataFrame () constructor

Parameter (type)	Description
DataFrame	The indices already in the old DataFrame are used, unless different ones are explicitly provided.

Pandas - Data structures

005-Pandas-Intro.ipynb

Index

- `Index` objects are responsible for storing axis labels and names.
- Any other array or sequence type provided as an index when building a `Series` or `DataFrame` object is internally converted to an `Index`.
- `Index` objects are immutable. This guarantees referential integrity when shared by different structures.
- `Index` provides methods and attributes to support set logic and value inspection.

Pandas - Data structures

Index subclasses



Class	Description
<code>Index</code>	Array of generic Python objects.
<code>Int64Index</code>	Integers.
<code>MultiIndex</code>	Hierarchical index, representing multiple indexing levels in a single axis. Similar to a tuple array.
<code>DatetimeIndex</code>	Timestamp with nanosecond resolution.
<code>PeriodIndex</code>	Time periods.

Pandas - Data structures



Index methods

Method	Description
<code>append</code>	Concatenates additional <code>Index</code> objects, producing a new object.
<code>diff</code>	Set difference.
<code>intersection</code>	Set intersection.
<code>union</code>	Set union.
<code>isin</code>	Computes a boolean array marking whether each of the entries in an <code>Index</code> is included in another collection.
<code>delete</code>	Creates a new <code>Index</code> object by removing an element from the original one.

Pandas - Data structures



Index methods

Method	Description
<code>drop</code>	Creates a new <code>Index</code> object by removing a set of elements from the original one.
<code>insert</code>	Creates a new <code>Index</code> object by inserting a new element into the original one.
<code>is_monotonic</code>	True if each element is greater or equal than the previous one. Alias for <code>is_monotonic_increasing()</code> . An <code>is_monotonic_decreasing()</code> method is also provided.
<code>is_unique</code>	True if the index has no duplicate elements.
<code>unique</code>	Computes an array containing the unique elements in the index.

Pandas - Essential functions

005-Pandas-Intro.ipynb

- Reindexing: `reindex()` method.
- Removing entries: `drop()` method.
- Indexing, selection, and filtering: operator `[]` and attributes `DataFrame.loc` and `DataFrame.iloc`.
- Arithmetic operations: operators `+`, `-`, `*`, and `/`, and methods `add()`, `sub()`, `mul()` y `div()`.
- Functional application and mapping: basic methods (`mean()`, `sum()`, ...), method `Series.map()`, methods `DataFrame.apply()`, and `DataFrame.applymap()`.
- Sorting and classification: methods `sort_values()`, `sort_index()`, and `rank()`.
- Managing indices with duplicates: method `Index.is_unique()`.

Pandas - Essential functions



`reindex()` parameters

Parameter	Description
<code>index</code>	New sequence to use as index.
<code>method</code>	Interpolation method (<code>ffill</code> or <code>bfill</code>).
<code>fill_value</code>	Value to use as a placeholder for null data.
<code>limit</code>	Maximum number of elements to fill using interpolation.
<code>level</code>	Hierarchical level to reindex.
<code>copy</code>	Marks whether data should be copied in case the new and the old indices are equivalent (<code>True</code> by default).

Pandas - Essential functions



DataFrame indexing

Syntax	Description
<code>obj[val]</code>	Selects a column or subset of columns, except if <code>val</code> is an array or the <code>DataFrame</code> is boolean, in which case it filters columns.
<code>obj.loc[val]</code>	Selects a row or subset of rows by label.
<code>obj.loc[:, val]</code>	Selects a column or subset of columns, by label.
<code>obj.loc[val1, val2]</code>	Selects both rows and columns, by label.

Pandas - Essential functions



DataFrame indexing

Syntax	Description
<code>reindex()</code>	Reorganizes one or more axes according to new indices.
<code>xs()</code>	Returns a cross-section of the dataframe attending to its labels.
<code>icol()</code> / <code>irow()</code>	Selects a single row / column attending to its location.
<code>get_value()</code> / <code>set_value()</code>	Selects a single value attending to row/column labels.

Pandas - Descriptive statistics

Reduction parameters



Parameter	Description
<code>axis</code>	The axis on which to perform the reduction (<code>rows=0</code> , <code>columns=1</code>).
<code>skipna</code>	Whether to exclude null values. <code>True</code> by default.
<code>level</code>	Reduction grouping by <code>level</code> in hierarchical indices.

Pandas - Descriptive statistics

Methods



Method	Description
<code>count</code>	Number of non-null values.
<code>describe</code>	Computes several statistics of a <code>Series</code> or the columns of a <code>DataFrame</code> .
<code>min</code> / <code>max</code>	Minimum / maximum value.
<code>argmin</code> / <code>argmax</code>	Location of the minimum / maximum value.
<code>idxmin</code> / <code>idxmax</code>	Index (label) of the minimum / maximum value.
<code>quantile</code>	Returns the specified p-quantile.
<code>sum</code>	Summation.

Pandas - Descriptive statistics Methods



Method	Description
<code>mean</code>	Arithmetic mean.
<code>median</code>	Median.
<code>mad</code>	Mean absolute deviation.
<code>var</code>	Variance.
<code>std</code>	Standard deviation.
<code>skew</code>	Skewness (third standardized moment).
<code>kurt</code>	Kurtosis (fourth standardized moment).

Pandas - Descriptive statistics Methods



Method	Description
<code>cumsum</code>	Cumulative sum.
<code>cummin,</code> <code>cummax</code>	Cumulative minimum / maximum.
<code>cumprod</code>	Cumulative product.
<code>diff</code>	First order differences.
<code>pct_change</code>	Percentage change.

Pandas - Unknown data

005-Pandas-Intro.ipynb

- Unknown or null data are common in most data analysis applications.
- Pandas simplifies the management of null data, e.g., all descriptive statistics functions automatically omit unknown values.
- `numpy.nan` is used as the default placeholder for unknown data. `None` will also be treated as a null value by Pandas.

Pandas - Unknown data

`fillna()` parameters



Parameter	Description
<code>value</code>	Escalar value or dictionary to use for filling.
<code>method</code>	Interpolation type (" <code>ffill</code> " or " <code>bfill</code> ").
<code>axis</code>	Axis to fill (by default 0, i.e., rows).
<code>inplace</code>	Modifies the object inplace, instead of creating a new copy.
<code>limit</code>	Maximum number of consecutive values to fill using interpolation.

Pandas - Hierarchical indices

005-Pandas-Intro.ipynb

- Hierarchical indexing allows to works (at a conceptual level) using `DataFrame` objects with more than 2 dimensions.
- It is implemented through the addition of different *levels* to the row and/or column indices of the table.
- Hierarchical indices are implemented by the `MultiIndex` class.
- It is possible to apply reduction operations to the different levels in the index hierarchy, obtaining a `DataFrame` (instead of a `Series`) after the reduction.

Pandas - Other considerations

Integer indexing

- Pandas objects with integer indices can be confusing, as the semantics of position- and label-based indexation vary.

```
obj = Series( range(5), index=['a', 'b', 'c', 'd',  
    'e'] )
```

```
obj[-1]
```

??

```
obj = Series( range(5), index=range(5) )
```

```
obj[-1]
```

??

Pandas - Other considerations

Integer indexing

- For objects with an integer index, Pandas cannot decide whether the user wants to apply location- or label-based indexing.
- It solves the ambiguity by always using label-based indexing.
- Provides methods `Series.iat_value()` and `DataFrame.irow()` for positional-based indexing in integer-indexed data.

Pandas: Data I/O

Reading / writing text

006-Pandas-DataIO.ipynb

- Pandas provides a set of methods to create `DataFrame` objects from tabular data stored in text format. The most useful ones are `read_csv()` and `read_table()` (which are nowadays mostly equivalent), which include options for:
 - Using one or more columns in the text file as `DataFrame` indices.
 - Naming columns using parameters or extracting names from the file.
 - Inferring types and performing data conversions.
 - Parsing dates, including combining several columns into a single one.
 - Iterating over chunks of large files (to fit data into memory).
 - Cleaning data: ignoring some rows or columns, comments, etc.
- Automatic type inference implies that it is not required to specify column types. Management of data and other non-basic types takes extra effort.

Pandas: Data I/O



`read_csv / read_table` parameters

Parameter	Description
<code>path</code>	URL of the file to open.
<code>sep / delimiter</code>	Regular expression to use for separating fields.
<code>header</code>	Row number containing the column names (0 by default), or <code>None</code> if correlative integers should be used.
<code>index_col</code>	Columns to build the <code>DataFrame</code> index.
<code>names</code>	Names for the <code>DataFrame</code> columns. Combine with <code>header=None</code> .
<code>skiprows</code>	Number of rows to ignore at the beginning of the file, or list of row numbers to ignore.

Pandas: Data I/O



`read_csv / read_table` parameters

Parameter	Description
<code>na_values</code>	Values that should be considered marks of unknown data.
<code>comment</code>	Regular expression to mark the beginning of a comment.
<code>parse_dates</code>	Tries to parse dates to a <code>datetime</code> object. <code>False</code> by default. If <code>True</code> , Pandas tries to parse all columns as dates. Alternatively, it can be a list of specific columns to parse as dates. If an element in the list is a tuple/list, it will try to combine the specified columns to parse a single date.
<code>keep_date_col</code>	If several columns are used to build a single date column, do not keep the joined columns. <code>True</code> by default.

Pandas: Data I/O



`read_csv / read_table` parameters

Parameter	Description
<code>converters</code>	Dictionary containing a mapping of columns to parsing functions. Each corresponding function will be applied to all the elements of a given column and the result will be inserted into the <code>DataFrame</code> .
<code>dayfirst</code>	When parsing potentially ambiguous dates, assume dates in international format (DD/MM/YYYY). <code>False</code> by default.
<code>date_parser</code>	Function to use to parse dates.
<code>nrows</code>	Number of rows to read, starting from the beginning of the file.
<code>iterator</code>	If <code>True</code> , the return value will be a <code>TextFileReader</code> object, to process the file in a chunk-by-chunk manner.

Pandas: Data I/O



`read_csv / read_table` parameters

Parameter	Description
<code>chunksize</code>	If chunk-by-chunk processing is active, size of each chunk.
<code>skip_footer</code>	Number of rows to ignore at the end of the file.
<code>verbose</code>	If <code>True</code> , print information about the parsing process.
<code>encoding</code>	Character encoding, e.g., 'utf-8'.
<code>squeeze</code>	If <code>True</code> , in case the result contains a single column, return a <code>Series</code> .
<code>thousands</code>	Thousands separator, e.g., ',' or '.'

Pandas: Data I/O

Binary formats

006-Pandas-DataIO.ipynb

- The `pickle` module in the Python standard library provides a convenient method to serialize (marshall) and deserialize (unmarshall) objects to binary format.
- Pandas provides the methods `to_pickle()` / `read_pickle()` which store/read data to/from `pickle` files.
- Using `pickle` is discouraged for long term storage, since `pickle` does not guarantee backwards compatibility.
- Pandas also provides methods to read and write Excel, HDF5, Stata, and HTML files, among others.

Pandas: Data I/O

HTML and web services

006-Pandas-DataIO.ipynb

- Many web sites include APIs that provide access to data sources in JSON format.
- There are several ways to access these services from Python.
- A simple way is through the `requests` package.
- The responses to requests will be converted to a JSON object.
- It is trivial to build a Pandas object from JSON, as we have seen.

Pandas: Data I/O Databases

006-Pandas-DataIO.ipynb

- Text and binary files are ultimately inefficient to store large amounts of data.
- Databases, both relational and non-relational, are one of the most common data sources in computer science.
- Pandas provides methods to load data from SQL queries. The `pandas.io.sql` module allows to execute SQL sentences and process the results.
- Other types of non-SQL databases, such as MongoDB, store objects in different formats: JSON, text, etc. Loading data mechanisms vary on a per-case basis.

- Much of the programming work in data analysis is spent on data preparation: loading, cleaning, transforming, and rearranging.
- Pandas and Python provide a set of flexible, high level data manipulation tools suitable for these preparation tasks:
 - Merging and combination of data: `merge()`, `concat()`, `combine_first()`.
 - Reshaping and pivoting: `stack()`, `unstack()`, `pivot()`.
 - Data transformation: removing duplicates, functional application, replacing values, renaming axes, discretization and binning, detecting and filtering outliers, permutation, and random sampling.
 - String manipulation: normalizing, cleaning, use of regular expressions...

Pandas: Data wrangling

Combining and merging datasets



007-Pandas-Data-wrangling.ipynb

- `pandas.merge()` connects rows in `DataFrame` based on one or more keys. It is similar to a `join` operation in a relational database.
- `pandas.concat()` glues or stacks together objects along an axis:
 - How are the new axes labeled? (union, intersection, ...)
 - Are the original groups identifiable in the resulting object?
 - Which axis should be used for concatenating?
- `combine_first()` enables splicing together overlapping data to fill in missing values in one object with values from the other.

Pandas: Data wrangling

`merge()` parameters



Parameter	Description
<code>left</code>	Left hand side operand of the merge.
<code>right</code>	Right hand side operand of the merge.
<code>how</code>	<code>['inner' 'outer' 'left' 'right']</code> . <code>'inner'</code> by default.
<code>on</code>	Column names to join on. Must be found in both <code>DataFrame</code> objects. If not specified, will use all columns with matching names.

Pandas: Data wrangling

`merge()` parameters



Parameter	Description
<code>left_on</code>	Columns in <i>left</i> operand to use as join keys.
<code>right_on</code>	Columns in <i>right</i> operand to use as join keys.
<code>left_index</code>	<code>[True False]</code> . Use row index in <i>left</i> as its join key.
<code>right_index</code>	<code>[True False]</code> . Use row index in <i>right</i> as its join key.

Pandas: Data wrangling

`merge()` parameters



Parameter	Description
<code>sort</code>	<code>[True False]</code> . Sort merged data lexicographically by join keys. <code>True</code> by default. Disable to get better performance in some cases on large datasets.
<code>suffixes</code>	Tuple of string values to append to column names in case of overlap. <code>('_x' , '_y')</code> by default.
<code>copy</code>	If <code>False</code> , avoid copying data into resulting data structure in some exceptional cases. By default always copies.

Pandas: Data wrangling

`concat()` parameters



Parameter	Description
<code>objs</code>	List or dictionary of Pandas objects to be concatenated. The only required argument.
<code>axis</code>	Axis to concatenate along. 0 by default (rows).
<code>join</code>	[<code>'inner'</code> <code>'outer'</code>]. <code>'outer'</code> by default. Whether to intersect (<code>inner</code>) or union (<code>outer</code>) together indices along the other axes.
<code>join_axes</code>	Specific indices to use for the other $(n-1)$ axes instead of performing union/intersection logic.

Pandas: Data wrangling


`concat()` parameters



Parameter	Description
<code>keys</code>	Values to associate with objects being concatenated, forming a hierarchical index along the concatenation axis.
<code>levels</code>	Specific indices to use as hierarchical index level or levels if <code>keys</code> passed.
<code>names</code>	Names for created hierarchical levels if <code>keys</code> and/or <code>levels</code> passed.
<code>verify_integrity</code>	Check new axis in concatenated objects for duplicates and raise exception if so. <code>False</code> by default.
<code>ignore_index</code>	Do not preserve indices along concatenation axis.

Pandas: Data wrangling

Reshaping and pivoting




007-Pandas-Data-wrangling.ipynb

- A different kind of transformations is its reshaping, sometimes called pivoting.
- It consists in transposing rows and columns, modifying data dimensionality.
- The basic reshaping and pivoting operations provided by Pandas are:
 - `stack()`: pivots from the columns in the data to the rows.
 - `unstack()`: pivots from the rows into the columns.
 - `pivot()`: reshapes rows and columns, allowing to transform tables in “long” format to “wide” format in a single step.

- We have focused on structural modifications. There are other types of transformations which focus on the data:
 - Removing duplicates: `duplicated()`, `drop_duplicates()`.
 - Functional application: `apply()`, `map()`, `applymap()`.
 - Value substitution: `replace()`.
 - Index renaming: `rename()`.
 - Discretization and binning: `cut()`, `qcut()`.
 - Detecting and filtering outliers.
 - Permutation and random sampling: `permutation()`, `take()`.
 - Computing indicators/dummy variables: `get_dummies()`.

Pandas: Data wrangling

String manipulation



007-Pandas-Data-wrangling.ipynb

- One of Python's most popular characteristics is its string manipulation routines.
- The `str` class provides methods for conveniently performing many string operations, such as searching, substitutions, splitting, etc.
- Many of these operations accept regular expressions as parameters.
- Pandas adds string functionality, as it allows to apply string operations over data tables automatically handling unknown values.

Method	Description
<code>count</code>	Returns the number of non-overlapping occurrences of substring in the string.
<code>endswith / startswith</code>	Returns <code>True</code> if a string ends with suffix / starts with prefix.
<code>join</code>	Use string as a delimiter for concatenating a sequence of other strings.
<code>index</code>	Return position of first character in substring if found in the string. Raises <code>ValueError</code> if not found.

Pandas: Data wrangling

str methods



Method	Description
<code>find</code>	Like <code>index()</code> , but returns <code>-1</code> if not found.
<code>rfind</code>	Like <code>find()</code> , but returns position of last occurrence.
<code>replace</code>	Replace occurrences of string with another string.
<code>strip</code> / <code>rstrip</code> / <code>lstrip</code>	Trim whitespace, including newlines.
<code>split</code>	Break string into list of substrings using passed delimiter.

Pandas: Data wrangling

str methods



Method	Description
<code>lower</code> / <code>upper</code>	Convert alphabet characters to lowercase or uppercase, respectively.
<code>ljust</code> / <code>rjust</code>	Left / right justify. Pad opposite side of string with spaces (or some other fill character) to return a string with a minimum width.

Pandas: Data wrangling

Regular expressions



007-Pandas-Data-wrangling.ipynb

- Regular expressions provide a flexible way to search or match string patterns in text.
- A regular expression is a string containing certain elements, called “special characters”, with a specific semantic.
- A regular expression (*regex*) describes a pattern to match in the text, or a way to manipulate it.
- The `re` module includes functions for regular expression application. It has three types of functions: pattern matching, substitutions, and splitting.

Pandas: Data wrangling regex methods



Method	Description
<code>findall / finditer</code>	Return all non-overlapping matching patterns in a string as a list / iterator.
<code>match</code>	Match pattern at start of string and optionally segment pattern components into groups. Returns a <code>match</code> object, or <code>None</code> .
<code>search</code>	Scan string for match to pattern; returning a <code>match</code> object if so. The match can be anywhere in the string, as opposed to <code>match()</code> .
<code>split</code>	Break string into pieces at each occurrence of pattern.
<code>sub, subn</code>	Replace all / first n occurrences of pattern with replacement expression.

Pandas: Data wrangling

Vectorized string functions in Pandas



Method	Description
<code>cat</code>	Concatenate strings element-wise with optional delimiter.
<code>contains</code>	Return boolean array if each string contains pattern.
<code>count</code>	Count occurrences of pattern.
<code>endswith</code> / <code>startswith</code>	Applies homonym functions in <code>str</code> element-wise.
<code>findall</code>	Compute list of all occurrences of pattern for each string.

Pandas: Data wrangling

Vectorized string functions in Pandas



Method	Description
<code>get</code>	Index into each element (retrieve i-th element).
<code>join</code>	Join strings in each element of the Series with passed separator.
<code>len</code>	Compute length of each string.
<code>lower / upper</code>	Convert cases.
<code>match</code>	Use <code>re.match</code> with passed regex on each element.
<code>pad</code>	Add whitespace to left, right, or both sides of string.

Pandas: Data wrangling

Vectorized string functions in Pandas



Method	Description
<code>center</code>	Equivalent to <code>pad(side="both")</code> .
<code>repeat</code>	Duplicate values.
<code>replace</code>	Replace occurrences of pattern with some other regex.
<code>slice</code>	Slice each string in the Series.
<code>split</code>	Split strings on delimiter.
<code>strip</code> / <code>rstrip</code> / <code>lstrip</code>	Trim whitespace, including newlines, element-wise.

Pandas: Plotting and visualization



008-Pandas-Plotting.ipynb

- Matplotlib provides a powerful framework for plotting and visualization, but it is a low-level tool.
- Building a plot involves configuring several Python objects without implicit semantics.
- Pandas objects are a centralized storage of data, with semantics at least partially known.
- Pandas objects provide a `plot()` method which builds complex plots in a convenient manner.

Pandas: Plotting and visualization



`plot()` parameters

Parameter	Description
<code>ax</code>	Matplotlib subplot object to plot on. If nothing passed, uses active subplot.
<code>kind</code>	[<code>'line'</code> <code>'bar'</code> <code>'barh'</code> <code>'hist'</code> <code>'box'</code> <code>'kde'</code> <code>'density'</code> <code>'area'</code> <code>'pie'</code> <code>'scatter'</code> <code>'hexbin'</code>]. Type of chart.
<code>logx / logy</code>	Use logarithmic scale on the X / Y axis.
<code>use_index</code>	Use the object index for tick labels.
<code>rot</code>	Rotation angle of tick labels.
<code>xticks / yticks</code>	Values to use for X / Y axis ticks.

Pandas: Plotting and visualization



`plot()` parameters

Parameter	Description
<code>xlim / ylim</code>	X / Y axis limits (specified as <code>[min, max]</code>).
<code>grid</code>	<code>[True False]</code> . Display axis grid.
<code>subplots</code>	<code>[True False]</code> . Plot each DataFrame column in a separate subplot.
<code>sharex / sharey</code>	<code>[True False]</code> . If <code>subplots=True</code> , share the same X / Y axis, linking ticks and limits.
<code>layout</code>	Tuple indicating the geometry of subplots to use.
<code>figsize</code>	Figure size.

Pandas: Plotting and visualization



`plot()` parameters

Parameter	Description
<code>title</code>	Plot title.
<code>legend</code>	<code>[True False 'reverse']</code> . Add a subplot legend.
<code>style</code>	Dictionary matching each column with the style to use for plotting it.
<code>loglog</code>	Use log scale for both axes.
<code>fontsize</code>	Font size to use for ticks.
<code>colormap</code>	Color map to index.

Pandas: Plotting and visualization



`plot()` parameters

Parameter	Description
<code>colorbar</code>	<code>[True False]</code> . Whether to draw the value legend for 'scatter' and 'hexbin' plot types).
<code>table</code>	If a Series or DataFrame object is provided, includes it in the plot. Useful to combine plots and tables.
<code>stacked</code>	Create stacked plot.
<code>sort_columns</code>	<code>[True False]</code> . Whether to sort columns lexicographically before plotting.
<code>secondary_y</code>	List of columns that should be referenced to a secondary Y axis.

Pandas: Plotting and visualization



`plot()` parameters

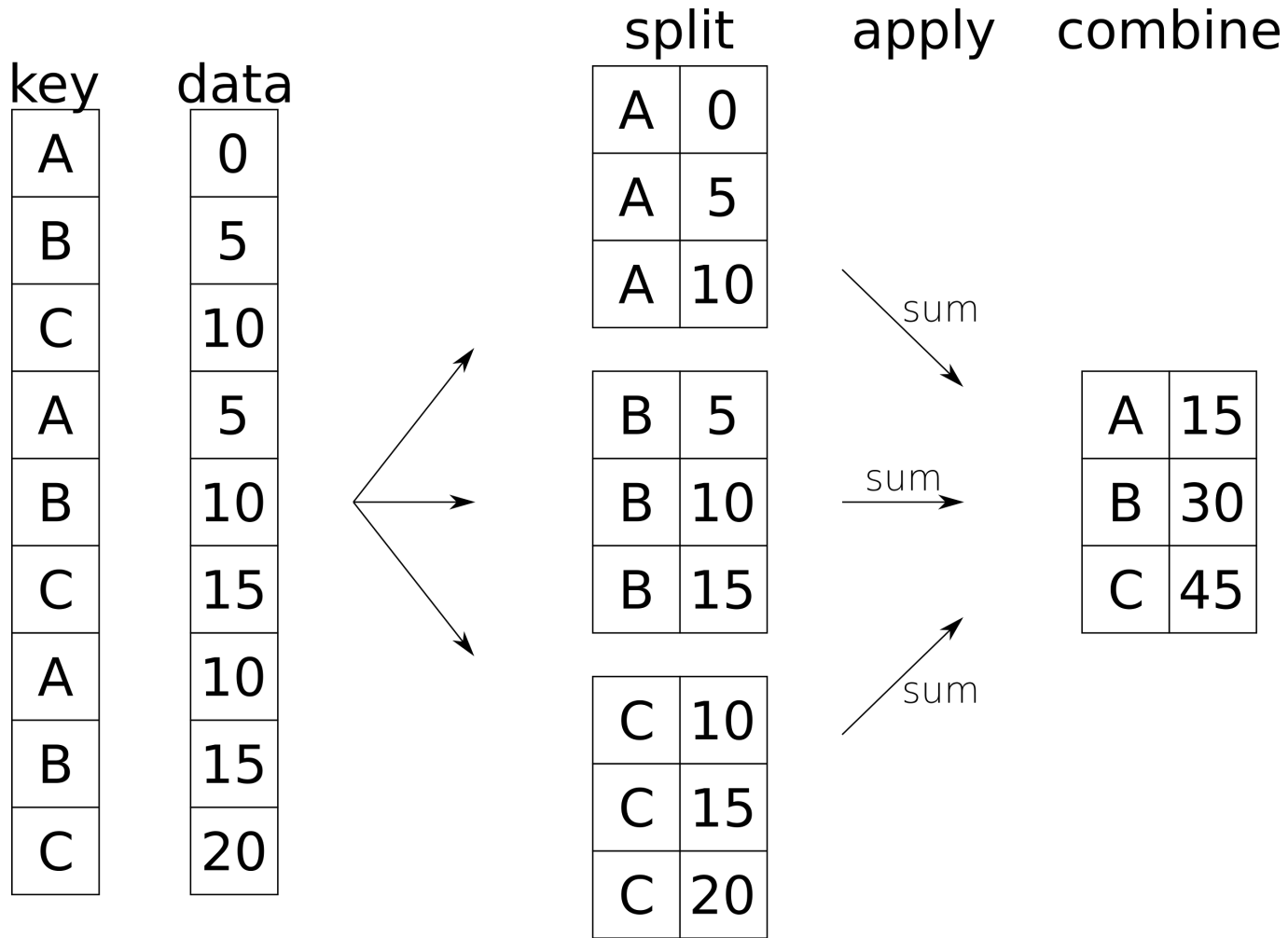
Parameter	Description
<code>mark_right</code>	<code>[True False]</code> . When a secondary axis is used, automatically add the suffix <code>"(right)"</code> to the legends of the series referenced to it.
<code>**kwargs</code>	Parameters not processed by Pandas will be passed to Matplotlib.

Pandas: Data aggregation



- One of the reasons for the popularity of relational databases and SQL is the ease with which data can be joined, filtered, transformed, and aggregated.
- However, query languages like SQL have limited expressiveness. Pandas allows to implement *split-apply-combine* operations conveniently:
 - Split a Pandas object into pieces using one or more keys.
 - Compute group summary statistics.
 - Apply a varying set of functions to each column of a `DataFrame`.
 - Compute pivot tables and cross-tabulations.
 - Perform quantile analysis and other data-derived group analyses.

Pandas: Data aggregation



Pandas: Data aggregation



009-Pandas-Aggregation.ipynb

- A critical aspect of these transformations is how to categorize data.
- A Pandas object is *split* into groups based on one or more keys, applied on a particular axis.
- The Pandas mechanism for this operation is `groupby()`.
- A grouping key can take many forms, e.g.:
 - A list or array of values that is the same length as the axis being grouped.
 - A value indicating a column name in a `DataFrame`.
 - A `dict` or `Series` giving correspondences between the values on the axis being grouped and the group names.
 - A function to be invoked on the axis index or the individual labels of the index.

Pandas: Data aggregation

Aggregation functions



009-Pandas-Aggregation.ipynb

- An *aggregation functions* is any transformation which produces a scalar value from an array (also called *reduction*).
- Aggregation functions as implemented by the `GroupBy` class have been optimized and are computed on the original data of the `DataFrame` or `Series`.
- Applicable aggregation functions are not limited to this subset: any function, including user-defined functions, can be applied to a grouped dataset.

Pandas: Data aggregation

GroupBy methods



Method	Description
count	Number of non-NA values in the group.
sum	Sum of non-NA values.
mean	Mean of non-NA values.
median	Arithmetic median of non-NA values.

Pandas: Agregación

Métodos en GroupBy



Método	Descripción
<code>std / var</code>	Unbiased standard deviation / variance.
<code>min / max</code>	Minimum / maximum of non-NA values.
<code>prod</code>	Product of non-NA values.
<code>first / last</code>	First / last non-NA value.

Pandas: Group-wise operations



009-Pandas-Aggregation.ipynb

- Aggregation is only one kind of group operation: accepts functions that reduce a one-dimensional array to a scalar value.
- In the general case, we want to apply any kind of operation to grouped data.
- This is done using `transform()` and `apply()`:
 - `transform()` broadcasts the result of an aggregation over the original members of the group.
 - `apply()` applies a function to each group and combines the results using `pandas.concat()`.

Pandas: Pivot tables and cross-tabulation



- A pivot table is a data summarization tool which aggregates a table by one or more keys, arranging the data in a rectangle with some groups along rows and some along columns.
- It can be built using `groupby()`, but `pivot_table()` provides a more convenient high-level interface.
- A cross-tabulation is a special case of a pivot table that computes group frequencies.
- Could also be built manually using several functions, but `crosstab()` simplifies the process.

Pandas: Time series



010-Pandas-Time-series.ipynb

- Any dataset which includes observations at many points in time forms a time series.
- Many time series are *fixed frequency*: data points occur at regular intervals.
- Others are *irregular*: without a fixed offset between data points.
- How time series are referred depends on the application. Among others:
 - *Timestamps*: specific instants in time.
 - *Fixed periods*, such as the month of January 2007 or the full year 2010.
 - *Intervals* of time, indicated by start and end timestamps.
 - Elapsed time relative to a particular fixed start time.
- Pandas provides a standard set of time series tools and data algorithms to slice and dice, aggregate, resample, etc.

Pandas: Time series

Base frequencies



Alias	Offset type	Description
D	Day	Calendar daily.
B	BusinessDay	Business daily.
H	Hour	Hourly.
T / min	Minute	Minutely.
S	Second	Secondly.
L / ms	Milli	Millisecond.
U	Micro	Microsecond.

Pandas: Time series

Base frequencies



Alias	Offset type	Description
M	MonthEnd	Last calendar day of month.
BM	BusinessMonthEnd	Last business day of month.
MS	MonthBegin	First calendar day of month.
BMS	BusinessMonthBegin	First business day of month.
W-MON, W-TUE, ...	Week	Weekly on given day of week: MON, TUE, WED, THU, FRI, SAT, or SUN.
WOM-1MON , WOM-1TUE, ...	WeekOfMonth	Generate weekly dates on the first, second, third, or fourth week of the month. For example, WOM-3FRI for the 3rd Friday of each month.

Pandas: Time series

Base frequencies



Alias	Offset type	Description
Q-JAN, ...	QuarterEnd	Quarterly dates anchored on last calendar day of each month, for year ending in indicated month.
BQ-JAN, ...	BusinessQuarterEnd	Quarterly dates anchored on last business day of each month, for year ending in indicated month.
QS-JAN, ...	QuarterBegin	Quarterly dates anchored on first calendar day of each month, for year ending in indicated month.
BQS-JAN, ...	BusinessQuarterBegin	Quarterly dates anchored on first business day of each month, for year ending in indicated month.

Pandas: Time series

Base frequencies



Alias	Offset type	Description
A-JAN, ...	YearEnd	Annual dates anchored on last calendar day of given month.
BA-JAN, ...	BusinessYearEnd	Annual dates anchored on last business day of given month.
AS-JAN, ...	YearBegin	Annual dates anchored on first calendar day of given month.
BAS-JAN , ...	BusinessYearBegin	Annual dates anchored on first business day of given month.

Pandas: Time series

`resample()` parameters



Parameters	Description
<code>freq</code>	String or <code>DateOffset</code> indicated desired resample frequency.
<code>how</code>	Function name or array function producing aggregated value.
<code>axis</code>	Axis to resample on, default to 0 (rows).
<code>fill_method</code>	How to interpolate when upsampling (<code>'ffill'</code> or <code>'bfill'</code>).
<code>closed</code>	In downsampling, which end of each interval is closed (inclusive). Defaults to <code>'right'</code> .
<code>label</code>	In downsampling, how to label the aggregated result, with the right or left bin edge. Defaults to <code>'right'</code> .

Pandas: Time series



resample() parameters

Parameters	Description
<code>loffset</code>	Time adjustment to the bin labels, such as <code>'-1s'</code> / <code>Second(-1)</code> to shift the aggregate labels one second earlier.
<code>limit</code>	When forward or backward filling, the maximum number of periods to fill.
<code>kind</code>	Aggregate to periods (<code>'period'</code>) or timestamps (<code>'timestamp'</code>); defaults to kind of index the time series has.
<code>convention</code>	When resampling periods, the convention (<code>'start'</code> or <code>'end'</code>) for converting the low frequency period to high frequency. Defaults to <code>'end'</code> .

Pandas: Time series

Moving window functions



Función	Descripción
<code>rolling_count</code>	Returns number of non-NA observations in each trailing window.
<code>rolling_sum</code>	Moving window sum.
<code>rolling_mean</code>	Moving window mean.
<code>rolling_median</code>	Moving window median.
<code>rolling_std / rolling_var</code>	Moving window variance / standard deviation. Uses (n-1) denominator.
<code>rolling_skew / rolling_kurt</code>	Moving window skewness (3rd moment) / kurtosis (4th moment).

Pandas: Time series

Moving window functions



Function	Description
<code>rolling_min / rolling_max</code>	Moving window minimum / maximum.
<code>rolling_quant ile</code>	Moving window score at percentile / sample quantile.
<code>rolling_corr / rolling_cov</code>	Moving window correlation / covariance.
<code>rolling_apply</code>	Apply generic array function over a moving window.

Pandas: Time series

Moving window functions



Function	Description
<code>ewma</code>	Exponentially-weighted moving average.
<code>ewmstd</code> / <code>ewmvar</code>	Exponentially-weighted moving standard deviation / variance.
<code>ewmcorr</code> / <code>ewmcov</code>	Exponentially-weighted moving correlation / covariance.

Bibliography



- *Python for Data Analysis: Data Wrangling with Pandas, NumPy, and Python*. Wes McKinney. O'Reilly Media, 2017 (2ª edición).
- *Python Data Science Handbook: Essential Tools for Working with Data*. Jake VanderPlas. O'Reilly Media, 2016 (1ª edición).

High level tools

Contents

- Visualization: Seaborn.
- Machine learning: Scikit-Learn.
- Out-of-core computation: Dask.



Seaborn

- Seaborn is a graphics library that improves the design and appearance of Matplotlib plots and includes additional statistical methods:
 - It allows to use themes to unify plot aesthetics.
 - Adds functions to visualize and compare distributions on one and two variables.
 - Linear regression tools.
 - Functions for visualizing data matrices and clustering.
 - Plotting statistical time series.
- Similar to ggplot in R.

Scikit-Learn



- Extension to SciPy (**SciPy Toolkit**) focused on machine learning.
- Compatible with tabular data from other libraries:
 - NumPy
 - Pandas
 - SciPy.sparse (matrices dispersas).
- Includes algorithms for supervised / unsupervised learning: GLMs, SVMs, kNN, Bayes, decision trees, clustering, etc.
- Cross-validation.
- Grid search of optimal models.
- Parallelization.



Dask

- Limited but simple alternative to Spark.
- Restrictions to Big Data in Python:
 - Parallelization: restricted by the GIL.
 - Physical memory.
- Dask provides two main tools:
 - Dynamic task planifier (cluster management).
 - “Big data” collections: trivially perform out-of-core computing.
- Provides limited out-of-core support to other APIs:
 - Pandas
 - Scikit-Learn
- Allows to parallelize computations on a cluster.

Python: Other tools

- Statistical modeling:
 - Statsmodels.
- Big Data:
 - Hadoop.
 - Spark.
- Deep Learning:
 - TensorFlow.
 - PyTorch.
- GPU data science:
 - RAPIDS.

Bibliography



- Python Data Science Handbook: Essential Tools for Working with Data. Jake VanderPlas. O'Reilly Media, 2016 (1ª edición).
- Data Science with Python and Dask. Jesse C. Daniel. Manning Publications, 2019 (1ª edición).