

# Data Analytics with Python

Pau, France

25 - 28 Octobre, 2017



Co-funded by the  
Erasmus+ Programme  
of the European Union



# Python

- General-purpose language, interpreted, object-oriented.
- Portable: an interpreter exist 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 - Executing programs

## 1. Using the interactive interpreter:

```
>>> txt = "Luck, I am your father"  
>>> print txt  
Luck, I am your father
```



# Python - Executing programs

## 2. Interpreting Python code:

```
----- > script.py  
txt = "Luck, I am your father"  
print txt  
-----
```

```
$ python2 script.py  
Luck, I am your father
```



# Python - Executing programs

## 3. Shell script:

```
----- > script.py  
#!/usr/bin/python2  
txt = "Luck, I am your father"  
print txt  
-----
```

```
$ ./script.py  
Luck, I am your father
```



# Python - Executing programs

## 4. Embedding Python into C:

```
----- > prueba.c
#include <Python.h>
...
Py_Initialize();
PyRun_SimpleString( "txt = 'Luck, I am your father'" );
PyRun_SimpleString( "print txt" );
...
-----
$ gcc -I /usr/include/python2.7 prueba.c -lpython2.7
$ ./a.out
Luck, 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 = "Luck, 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:

```
----- > modulo1.py
```

```
PI = 3.14159
```

```
def f(x):  
    return x + 2
```

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

```
-----
```



# Python - Semantics & Syntax

## Binary operators

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



# Python - Scalar types

None	“null” in Python. Singleton class.
str	String type. ASCII in Python 2.x, Unicode in Python 3.
unicode	Unicode string.
float	Floating point number, double precision (64 bits).
bool	Logical, True or False.
int	Signed integer. 32 or 64 bits depending on platform.
long	Signed integer, arbitrary precision.



# Python - Scalar types

## Numerical types

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

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

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

```
>>> fval = a_int / float(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 0x7fff664159510>
>>> 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 &amp; 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 )  
17.9900000000000002
```



# 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 ( $\lambda$ ) 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 ( $\lambda$ ) 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)))
```

```
530
```



# 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
r	Read only.
w	Write only. Creates a new file, overwriting any previous one.
a	Concatenate to a file (created if it does not exist).
r+	Read-write.
b	Binary mode (usage example: <code>'rb'</code> ).
U	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>	True if the file is closed.



# iPython

- Alternative interactive Python interpreter.
- Does not provide any analysis or computation tool beyond those in cPython, but greatly simplifies some repetitive tools and the general development process.
- It includes GUI environments to simplify data visualization, e.g., it provides a web interface to Python, allowing to save HTML sessions which can then be shared (called Notebooks, we will use them extensively in this course).



# iPython

```
$ ipython2
```

```
Python 2.7.13 (default, Jul 21 2017, 03:24:34)
```

```
Type "copyright", "credits" or "license" for more information.
```

```
IPython 5.3.0 -- An enhanced Interactive Python.
```

```
? -> Introduction and overview of IPython's features.
```

```
%quickref -> Quick reference.
```

```
help -> Python's own help system.
```

```
object? -> Details about 'object', use 'object??' for extra details.
```

```
In [1]:
```



# iPython - Features

## Pretty printing

```
In [1]: a = 5
```

```
In [2]: a
```

```
Out[2]: 5
```

```
In[3]: from numpy.random import  
randn
```

```
In[4]: data = { i: randn() for i in  
range(7) }
```

```
In[4]: data
```

```
{0: -2.171772092191186,  
1: 0.2947739898361525,  
2: -0.3602586072105093,  
3: 1.300393521182511,  
4: -0.00033416357123924607,  
5: -0.572174661598875,  
6: 0.8141305611469151}
```



# iPython - Features

## Enumerated prompt

- iPython provides a prompt in which inputs and outputs are identified by an integer. This integer can be used to recover both:

```
In[5]: _i3
```

```
Out[5]: u'data = { i: randn() for i in range(7) }'
```

```
In[6]: _4
```

```
{0: -2.171772092191186,
```

```
 1: 0.2947739898361525,
```

```
...
```



# iPython - Features

## Autocompletion

- One of the improvements in usability w.r.t. cPython is the tab completion mechanism.
- When typing a command, pressing <Tab> activates a search of the current namespace, trying to identify any variable (objects, functions, etc.) whose name matches what has been typed so far:

```
In[1]: a_number = 1
```

```
In[2]: a_letter = 'a'
```

```
In[3]: a<Tab>
```

```
a_letter  a_number  abs          %alias      %alias_magic  
all       and       any          apply       as  
assert    %autocall %autoindent %automagic
```



# iPython - Features

## Autocompletion

- Besides finding the variables `a_letter` and `a_number`, in the previous example there are other valid completions, such as the magic commands `%alias` or `%autoindent`, as well as the operator `and` or the function `any()`.
- The autocompletion mechanism also allows to explore the methods and attributes of an object:

```
In[3]: b = [1, 2, 3]
```

```
In[4]: b.<Tab>
```

```
b.append      b.count      b.extend     b.index     b.insert  
b.pop         b.remove     b.reverse    b.sort
```



# iPython - Features

## Autocompletion

- It is also possible to use autocompletion with the scope of a module:

```
In[5]: import datetime
```

```
In[6]: datetime.<Tab>
```

```
datetime.MAXYEAR      datetime.MINYEAR      datetime.date  
datetime.datetime     datetime.datetime_CAPI datetime.time  
datetime.timedelta    datetime.tzinfo
```



# iPython - Features

## Autocompletion

- By default, methods and attributes starting with an underscore are hidden.
- To show them, we must include the underscore explicitly:

```
In[7]: datetime._<Tab>
```

```
datetime.__class__  
datetime.__file__  
datetime.__init__
```

```
datetime.__dict__  
datetime.__getattr__  
datetime.__new__ ...
```



# iPython - Features

## Autocompletion

- Besides looking for variables and functions, the autocompletion mechanism also searches file names when the interpreter deems it appropriate.
- It is also useful to explore keyword parameter names.



# iPython - Features

## Introspection

- Typing a question mark before or after a name returns information about it:

```
In[1]: b = [1,2,3]
```

```
In[2]: b?
```

```
Type:          list
```

```
String form:  [1, 2, 3]
```

```
Length:       3
```

```
Docstring:
```

```
list() -> new empty list
```

```
list(iterable) -> new list initialized from iterable's items
```



# iPython - Features

## Introspection

- If the name to be inspected is a function, it prints its docstring if available:

```
In[1]: import numpy as np
```

```
In[2]: np.linalg.det?
```

```
Type:          function
```

```
String form: <function det at 0x7f27c44b4ed8>
```

```
File:          /usr/lib/python2.7/site-packages/numpy/linalg/linalg.py
```

```
Definition:    numpy.linalg.det(a)
```

```
Docstring:
```

```
Compute the determinant of an array.
```

# iPython - Features

## Introspection

Parameters

-----

`a : (... , M, M) array_like`

Input array to compute determinants for.

Returns

-----

`det : (...) array_like`

Determinant of ``a``.

...

# iPython - Features

## Introspection

- If two question marks are used (e.g., `numpy.linalg.det??`), the source code of the function is shown, if available.
- The question mark character also allows to search a namespace using regular expressions:

```
In[3]: np.linalg.*rank*?  
numpy.linalg.matrix_rank
```



# iPython - Magic commands

- iPython defines many so-called “magic commands”, designed to simplify frequent tasks.
- The syntax of magic commands is `%magicname`.
- They behave like programs which are interpreted by iPython.
- Some magic commands receive arguments.

```
In[1]: %reset?
```

...

```
Resets the namespace by removing all names defined by the user
```

...



# iPython - Magic commands

`%run`

- `%run` executes Python code in a given file:

```
-----> test.py  
print "Hello world!"  
-----
```

```
In[1]: %run test.py  
Hello world!
```



# iPython - Magic commands

`%run`

- The file is executed inside an empty namespace. As such, the execution is equivalent to a command-line invocation.
- However, all the names defined in the script (imports, functions, global variables) will be accessible by the iPython environment after the `%run`.
- If the script expects arguments, this can be provided as usual:

```
In[1]: %run script.py <args>
```



# iPython - Magic commands

`%paste` / `%cpaste`

- A quick way to execute code is to copy-paste from a web page or other sources.
- Using `<Ctrl-C>` and `<Ctrl-V>` is tricky with Python, due to the special indentation semantics.:

```
x=5
```

```
y=7
```

```
if x > 5:
```

```
    x += 1
```

```
    y += 8
```

- The white line in the previous code prevents correct execution: some interpreters throw an `IndentationError`, while others execute `y += 8` outside the conditional block.



# iPython - Magic commands

`%paste` / `%cpaste`

- `%paste` and `%cpaste` allow to copy-paste code without worrying about indentation.
- `%paste` executes the text in the clipboard automatically.
- `%cpaste` is similar, but allows to edit the clipboard text before its execution.



# iPython - Magic commands

Command	Description
<code>%quickref</code>	Show a quick reference sheet.
<code>%magic</code>	Print information about the magic function system.
<code>%debug</code>	Activate the interactive debugger using the top of the stack from the last error.
<code>%hist</code>	Print input history.



# iPython - Magic commands

Command	Description
<code>%pdb</code>	Control the automatic calling of the <code>pdb</code> interactive debugger.
<code>%paste</code>	Paste & execute a pre-formatted code block from clipboard.
<code>%cpaste</code>	Opens a special prompt to copy and format code manually.
<code>%reset</code>	Resets the namespace by removing all names defined by the user.
<code>%page &lt;object&gt;</code>	Pretty print an object and display it through a pager.



# iPython - Magic commands

Command	Description
<code>%run script.py</code>	Executes the script in file <code>script.py</code> .
<code>%prun &lt;instruction&gt;</code>	Executes <code>&lt;instruction&gt;</code> using <code>cProfile</code> (Python profiler).
<code>%time &lt;instruction&gt;</code>	Measures the execution time of <code>&lt;instruction&gt;</code> .
<code>%timeit &lt;instruction&gt;</code>	Executes <code>&lt;instruction&gt;</code> multiple times and measures its execution time. Useful to time commands with short executions.



# iPython - Magic commands

Command	Description
<code>%who / %who_ls / %whos</code>	Print all interactive variables with different verbosity levels.
<code>%xdel &lt;variable&gt;</code>	Deletes <variable>.



# iPython - Shortcuts

Command	Description
<code>&lt;Ctrl+P&gt; / &lt;cursor up&gt;</code>	Search backwards in the command history, taking into account the current prompt.
<code>&lt;Ctrl+N&gt; / &lt;cursor down&gt;</code>	Search forward in the command history, taking into account the current prompt.
<code>&lt;Ctrl+R&gt;</code>	Search backwards in the command history in the readline style.
<code>&lt;Ctrl+Shift+V&gt;</code>	Paste text from the clipboard.
<code>&lt;Ctrl+C&gt;</code>	Interrupt current interactive execution.



# iPython - Shortcuts

Command	Description
<code>&lt;Ctrl+A&gt;</code>	Moves cursor to the beginning of the current line.
<code>&lt;Ctrl+E&gt;</code>	Moves cursor to the end of the current line.
<code>&lt;Ctrl+K&gt;</code>	Removes text from the cursor up to the current line end.
<code>&lt;Ctrl+U&gt;</code>	Removes all text from the current line.
<code>&lt;Ctrl+F&gt;</code> / <code>&lt;cursor derecho&gt;</code>	Move cursor forward one character.
<code>&lt;Ctrl+B&gt;</code> / <code>&lt;cursor izquierdo&gt;</code>	Move cursor backward one character.
<code>&lt;Ctrl+L&gt;</code>	Clear screen.



# iPython - Command history

## Search and reuse

- Using `<cursor up>` and `<cursor down>` it is possible to explore the command history. Contrary to the cPython interpreter, iPython takes into account the current prompt to restrict the search.
- Using `<Ctrl+R>` the current prompt is searched backwards in the command history. To discard a result and continue searching, press `<Ctrl+R>` again.

```
In[1]: a_command = f( x, y, z )
```

```
(reverse-i-search) 'com': a_command = f( x, y, z )
```



# iPython - Command history

## Input and output

- iPython stores both the input and output command history.
- The last two outputs are stored in variables `_` (an underscore) and `__` (two underscores), respectively.
- Input lines are stored in a set of variables named `_iX`, where `X` is the sequence number of the input. In a similar way, a set of variables named `_X` stores the corresponding outputs.
- **Warning:** these logs reference all the outputs of the interactive session, and so they will never be garbage collected. This may result in memory issues, particularly in long sessions, which must be manually solved through `%xdel` and/or `%reset`.



# iPython - Command history

## Input and output logs

- iPython allows to record an interactive session, including inputs and outputs.
- The log is activated using `%logstart`.
- It is deactivated using `%logoff`.
- The log can be temporarily paused and reactivated using `%logstop` and `%logon`, respectively.
- `%logstate` shows the current log state.



# iPython - Interacting with the OS

Command	Description
<code>!cmd</code>	Executes <code>cmd</code> in an OS shell.
<code>output = !cmd</code>	Executes <code>cmd</code> and stores its standard output into variable <code>output</code> .
<code>%alias alias_name cmd</code>	Defines an alias for command <code>cmd</code> with the name <code>alias_name</code> .
<code>%bookmark</code>	Folder bookmarking system in iPython.
<code>%cd &lt;dir&gt;</code>	Changes the working directory to <code>&lt;dir&gt;</code> .



# iPython - Interacting with the OS

Command	Description
<code>%pwd</code>	Prints the current working directory.
<code>%pushd &lt;dir&gt;</code>	Stacks the current working directory and changes it to <code>&lt;dir&gt;</code> .
<code>%popd</code>	Changes the current working directory to the top of the stack and pops it.
<code>%dirs</code>	Shows the current working directory stack.
<code>%dhist</code>	Shows the list of visited directories.
<code>%env</code>	Returns a dictionary with the environment.



# iPython - Software tools

## ipdb debugger

- The ipdb debugger included with iPython improves over the pdb standard by including autocompletion, syntax highlighting, and more verbose exception info.
- There are several ways to invoke ipdb:
  1. `%debug` enters the last error location.
  2. Activating the automatic debugging using `%pdb`.
  3. Executing `%debug <command>`.
  4. `%run -d script.py` launches the script through ipdb.
  5. Executing `import ipdb; ipdb.run( "<command>" )`
  6. Inserting a call to `ipdb.set_trace()` at any point in a Python script.



# iPython - Software tools

## ipdb commands

Command	Description
<code>h(elp)</code>	Shows the available commands.
<code>help &lt;command&gt;</code>	Shows the documentation available for <code>&lt;command&gt;</code> .
<code>c(ontinue)</code>	Resumes program execution.
<code>q(uit)</code>	Finishes the debugging session.
<code>b(reak) number</code>	Sets a breakpoint at line <code>number</code> of the current file.
<code>b path:number</code>	Sets a breakpoint at line <code>number</code> of file <code>path</code> .



# iPython - Software tools

## ipdb commands

Command	Description
<code>s (tep)</code>	Execute the current line, stop at the first possible occasion (either in a function that is called or in the current function).
<code>n (ext)</code>	Continue execution until the next line in the current function is reached or it returns.
<code>u (p) / d (own)</code>	Move the current frame one level up / down in the stack trace (to an older / newer frame).
<code>a (rgs)</code>	Print the arguments of the current function.



# iPython - Software tools

## ipdb commands

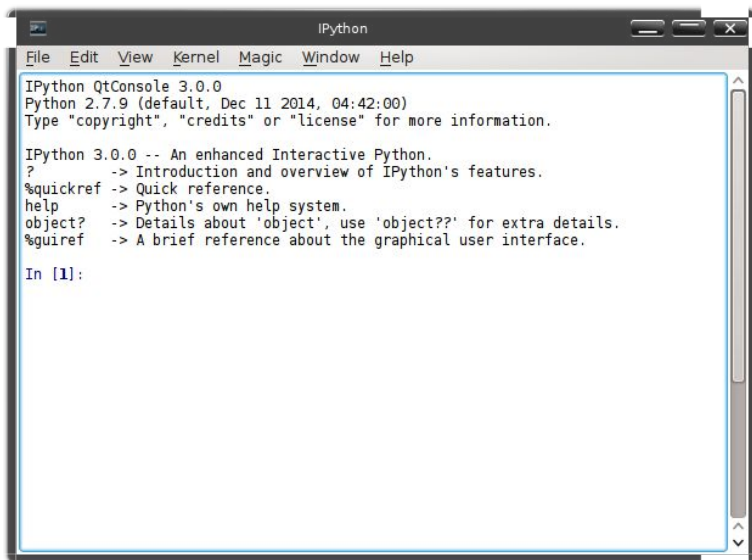
Command	Description
<code>debug instruction</code>	Enter a recursive debugger that steps through the code argument (which is an arbitrary expression or statement to be executed in the current environment).
<code>l(ist) [first [,last]]</code>	List source code for the current file, either around the current line or between lines <code>first</code> and <code>last</code> .
<code>w(here)</code>	Print the current stack trace.



# iPython - Software tools

## QT console

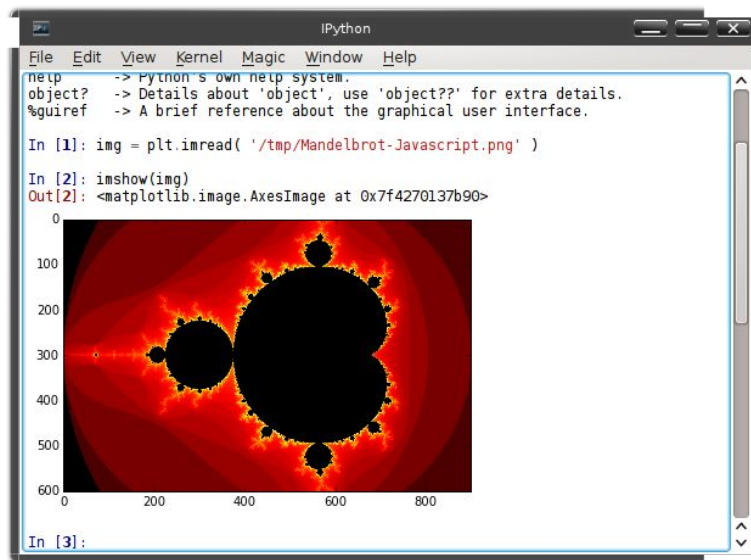
```
$ ipython2 qtconsole --pylab=inline
```



```
IPython QtConsole 3.0.0
Python 2.7.9 (default, Dec 11 2014, 04:42:00)
Type "copyright", "credits" or "license" for more information.

IPython 3.0.0 -- An enhanced Interactive Python.
?                -> Introduction and overview of IPython's features.
%quickref       -> Quick reference.
help            -> Python's own help system.
object?        -> Details about 'object', use 'object??' for extra details.
%gui?          -> A brief reference about the graphical user interface.

In [1]:
```



```
IPython QtConsole 3.0.0
Python 2.7.9 (default, Dec 11 2014, 04:42:00)
Type "copyright", "credits" or "license" for more information.

IPython 3.0.0 -- An enhanced Interactive Python.
?                -> Introduction and overview of IPython's features.
%quickref       -> Quick reference.
help            -> Python's own help system.
object?        -> Details about 'object', use 'object??' for extra details.
%gui?          -> A brief reference about the graphical user interface.

In [1]:
help         -> Python's own help system.
object?     -> Details about 'object', use 'object??' for extra details.
%gui?      -> A brief reference about the graphical user interface.

In [1]: img = plt.imread( '/tmp/Mandelbrot-Javascript.png' )
Out[2]: <matplotlib.image.AxesImage at 0x7f4270137b90>

In [2]: imshow(img)
Out[2]: <matplotlib.image.AxesImage at 0x7f4270137b90>

In [3]:
```



# iPython - Software tools

## QT console

- The QT console simplifies some tasks:
  - It allows to open several windows attached to a single iPython kernel.
  - Automatically shows the documentation of functions when interactively writing code.
- The `--pylab` option enables several environment changes. The underlying idea is to automatically get an analysis and development GUI similar to Matlab:
  - It loads `matplotlib` to draw plots.
  - Automatically imports `numpy`.
  - Enables support for embedded plots with `--inline`.



# iPython - Software tools

## HTML Notebook

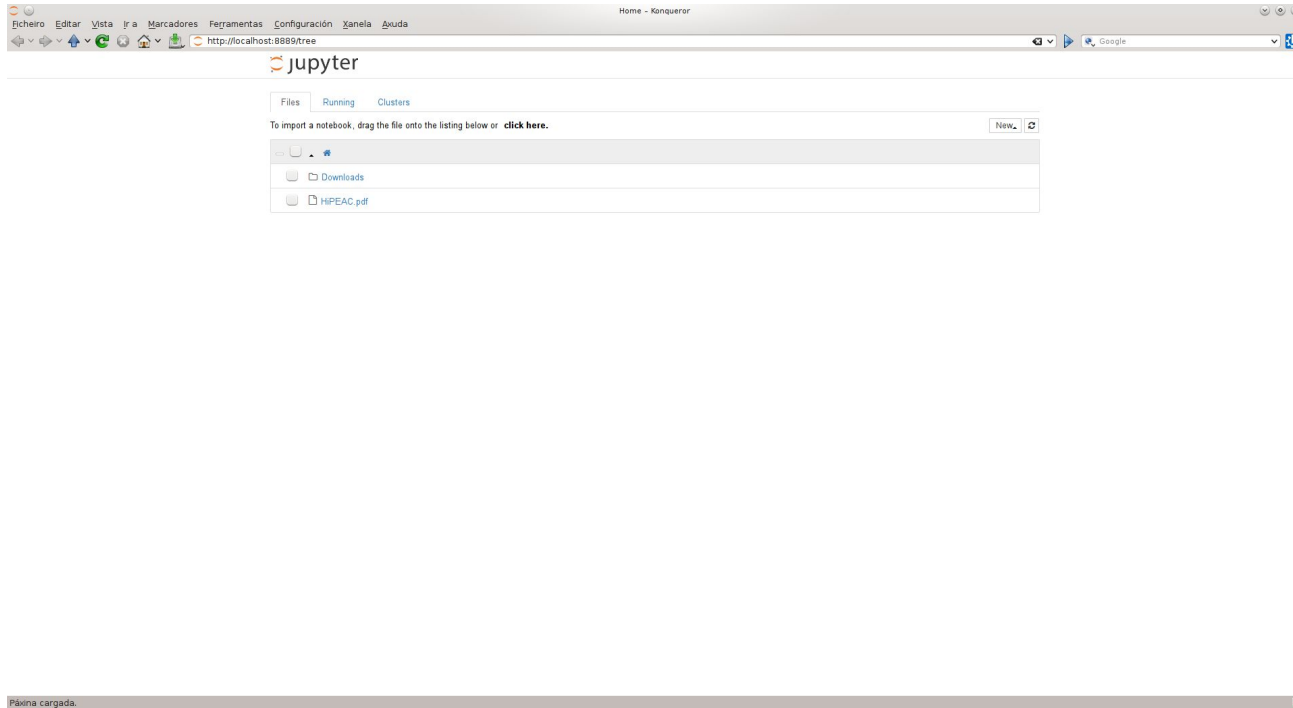
- iPython also includes an HTML environment known as Notebook.
- It uses a JSON-based format to share code, outputs and plots.
- It runs on a web server which can be accessed using a regular browser.
- Allows to use a remote iPython interpreter easily through the web.
- To launch a server/client pair locally, we run:

```
$ ipython2 notebook
```



# iPython - Software tools

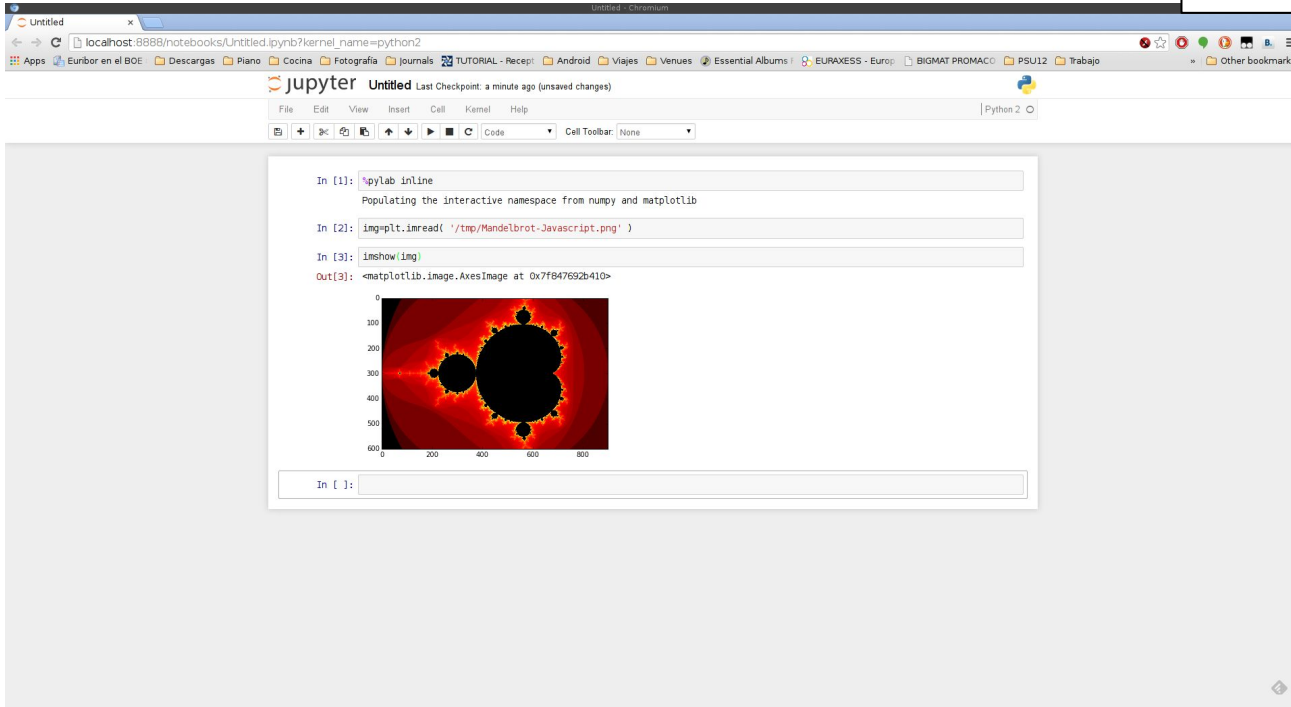
## HTML Notebook



# iPython - Software tools

## HTML Notebook

Notebook-001.ipynb



The screenshot shows a Jupyter Notebook running in a Chrome browser at localhost:8888. The notebook is titled "Untitled" and has a Python 2 kernel. The code in the notebook is as follows:

```
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 0x7f647692b410>
```

The output of the third cell is a plot of the Mandelbrot set, which is a fractal shape with a complex, self-similar boundary. The plot is displayed on a grid with x and y axes ranging from 0 to 600. The fractal is rendered in black and red, with the interior of the fractal being black and the exterior being red.



# 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

Notebook-002.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

Notebook-002.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

Notebook-002.ipynb

Función	Descripción
<pre>zeros() zeros_like()</pre>	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.
<pre>empty() empty_like()</pre>	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.
<pre>eye(), identity()</pre>	Build an array with 1s in its main diagonal and 0s elsewhere. <code>identity()</code> requires square dimensions.



# NumPy - ndarray Datatypes

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



# NumPy - ndarray Datatypes

Tipo	Código	Descripción
<code>float16</code>	<code>f2</code>	Floating point, half precision (2 bytes).
<code>float32</code>	<code>f4</code> o <code>f</code>	Floating point, single precision (4 bytes). Compatible with a C <code>float</code> .
<code>float64</code>	<code>f64</code> o <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> o <code>g</code>	Floating point, extended precision (16 bytes).



# NumPy - ndarray Datatypes

Tipo	Código	Descripción
<code>complex64</code>	<code>c8</code>	Complex number represented as 2 x <code>float32</code> .
<code>complex128</code>	<code>c16</code>	Complex number represented as 2 x <code>float64</code> .
<code>complex256</code>	<code>c32</code>	Complex number represented as 2 x <code>float128</code> .



# NumPy - ndarray Datatypes

Tipo	Código	Descripción
<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

Notebook-002.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

## Broadcasting

Notebook-002.ipynb

- 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

B : 3

A\*B: 256 x 256 x 3

A : 256 x 256 x 3

B : 256 x 1 x 3

A\*B: 256 x 256 x 3



# NumPy - ndarray

## Basic indexing

Notebook-002.ipynb

- 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

Notebook-002.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

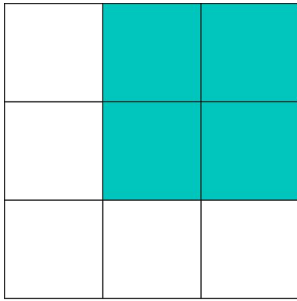


# NumPy - ndarray

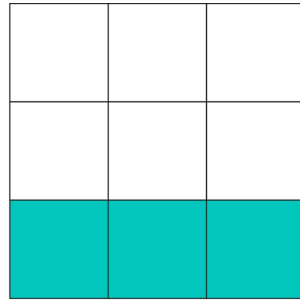
## Indexing with slices

Notebook-002.ipynb

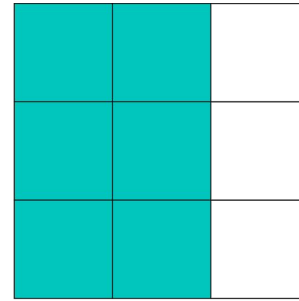
- The slice syntax is also valid for indexing arrays.
- Working with n-dimensional arrays, it is possible to use any combination of basic indexing and slicing.



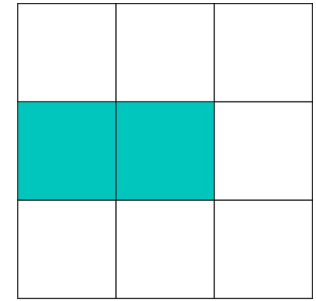
`arr[:2, 1:]`  
shape = (2,2)



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



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



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



# NumPy - ndarray

## Boolean indexing

Notebook-002.ipynb

- 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

## Fancy indexing

Notebook-002.ipynb

- 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 Transposition

Notebook-002.ipynb

- `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

Notebook-002.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

Notebook-002.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

Notebook-002.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

Notebook-002.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

Notebook-002.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

Notebook-002.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`)

Notebook-002.ipynb

Function	Description
<code>chisquare</code>	Samples a $\chi^2$ 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

Notebook-003.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

Notebook-003.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 $x$ rows and $y$ columns. Besides, it activates subfigure number $z$ ( $1 < z < x * y$ , row major order).



# Matplotlib - Introduction

Notebook-003.ipynb

Function	Description
<code>bar / barh</code>	Create vertical / horizontal bar charts. To plot stacked bars the <code>bottom</code> parameter 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

Notebook-003.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

Notebook-003.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

Notebook-003.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

Notebook-003.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

Notebook-003.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>	<a href="#">See docs</a>	Line type.
<code>linewidth</code>	<code>float</code>	Line width, in points.



# Matplotlib - Styles

Notebook-003.ipynb

Property	Value type	Description
<code>marker</code>	<a href="#">See docs</a>	Marker used for points in a line.
<code>mec</code>	Matplotlib color	Marker Edge Color.
<code>mew</code>	float	Marker Edge Width.
<code>mfc</code>	Matplotlib color	Marker Face Color.
<code>markersize</code>	float	Marker size, in points.



# Matplotlib - Styles

Notebook-003.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

Notebook-003.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

Notebook-003.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

Notebook-003.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

Notebook-003.ipynb

## Parameter to `hist()`

## Description

<code>histtype</code>	By default, bar chart. Other values: <ul style="list-style-type: none"><li>• <code>barstacked</code>: stacks bars when working with multiple data series.</li><li>• <code>step</code>: line without filling.</li><li>• <code>stepfilled</code>: line with filling.</li></ul>
<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

Notebook-003.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

Notebook-003.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.



# Text formats

## JSON

- JSON (*JavaScript Object Notation*) is an open standard and text format targetted at transmitting data in the form of key-value pairs.
- It is used as an alternative to HTML for transmitting data in client-server applications.
- Although originally derived from JavaScript, JSON is independent and many languages and interpreters include mechanisms to automate reading and writing JSON data.
- It is much more flexible than other text-based alternatives such as CSV.



# Text formats

## JSON

```
obj = ""  
{ "name": "John",  
  "residences": ["USA", "Spain", "France"],  
  "pet": null,  
  "relatives": [ {"name": "Andrew", "age": 25,  
                  "pet": "Zuko" },  
                  {"name": "Helen", "age": 33,  
                  "pet": "Cisco"} ]  
}  
""
```



# Text formats

## JSON

- Almost valid Python code.
- Basic types: objects (dictionaries), arrays (lists), strings, numbers, booleans, and nulls.
- Keys must be strings.
- The `json` module, included in the Python standard library, allows to read and write Python objects in JSON using its `loads()` and `dumps()` methods, respectively.



# Text formats

## HTML y XML

Notebook-004.ipynb

- Many web services communicate through HTML/XML.
- The `urllib2` module allows to open HTTP connections.
- The `lxml` allows to parse XML documents, e.g. HTML.
- Parsing non-semantic web pages requires detailed knowledge of their structure.
- Many HTML elements exist only for formatting purposes. In order to extract the underlying semantics it is necessary to process formats, syntax, etc. This process is called data wrangling, or data munging.



# 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

## Series

Notebook-005.ipynb

- 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

## DataFrame

Notebook-005.ipynb

- 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 DataFrame. 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 <code>2D 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

## Index

Notebook-005.ipynb

- `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

Notebook-005.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 (rows=0, columns=1).
<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 / max</code>	Minimum / maximum value.
<code>argmin / argmax</code>	Location of the minimum / maximum value.
<code>idxmin / 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
mean	Arithmetic mean.
median	Median.
mad	Mean absolute deviation.
var	Variance.
std	Standard deviation.
skew	Skewness (third standardized moment).
kurt	Kurtosis (fourth standardized moment).



# Pandas - Descriptive statistics Methods

Method	Description
<code>cumsum</code>	Cumulative sum.
<code>cummin, 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

Notebook-005.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

Notebook-005.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 indexing 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

Notebook-006.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

Notebook-006.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

Notebook-006.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

Notebook-006.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.



# Pandas: Data wrangling

- 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

Notebook-007.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

Notebook-007.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.



# Pandas: Data wrangling

## Data transformations

Notebook-007.ipynb

- 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

Notebook-007.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.



# Pandas: Data wrangling

## `str` methods

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>rstrip</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

Notebook-007.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 match 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 / 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>rstrip</code> / <code>lstrip</code>	Trim whitespace, including newlines, element-wise.



# Pandas: Plotting and visualization

Notebook-008.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'   'bar'   'barh'   'hist'   'box'   'kde'   'density'   'area'   'pie'   'scatter'   '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>	[True False 'reverse']. 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 <code>Series</code> or <code>DataFrame</code> 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>**kwds</code>	Parameters not processed by Pandas will be passed to Matplotlib.

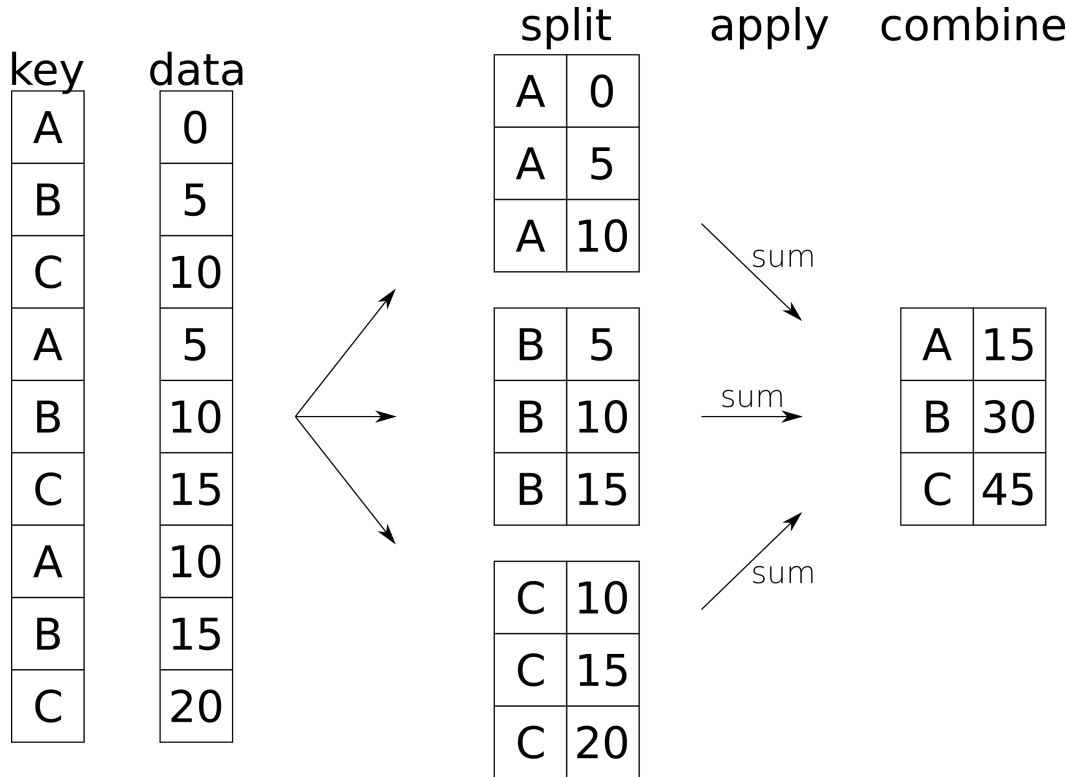


# Pandas: Data aggregation

- 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

Notebook-009.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

Notebook-009.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
<code>count</code>	Number of non-NA values in the group.
<code>sum</code>	Sum of non-NA values.
<code>mean</code>	Mean of non-NA values.
<code>median</code>	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

Notebook-009.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

Notebook-009.ipynb

- 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

Notebook-010.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

<b>Alias</b>	<b>Offset type</b>	<b>Description</b>
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

<b>Alias</b>	<b>Offset type</b>	<b>Description</b>
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
rolling_count	Returns number of non-NA observations in each trailing window.
rolling_sum	Moving window sum.
rolling_mean	Moving window mean.
rolling_median	Moving window meadian.
rolling_std / rolling_var	Moving window variance / standard deviation. Uses (n-1) denominator.
rolling_skew / rolling_kurt	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_quantile</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.



# Scientific data libraries

- A computer uses multiple I/O technologies:
  - Magnetic tapes for backups.
  - Magnetic disks for HDDs.
  - SSDs for fast access to system functions.
  - Flash (USB pen drives) for portability.
  - Optical disks (DVD) for software distribution.
- POSIX defines functions to access data independently from the **physical** storage medium (`read()`, `write()`, `fread()`, etc.)



# Scientific data libraries

- However, under POSIX the programmer is responsible for:
  - Using the OS-provided primitives (not all OS support POSIX!)
  - Organizing data according to a sequence-of-bytes model.
  - Performing the required conversions among different storage formats (e.g., when reading files generated by a *big-endian* architecture on a *little-endian* one).
- Basic aim: decouple applications from the storage **logic**.



# Scientific data libraries

## NetCDF

- In 1982, NASA designs CDF.
- Goal: to show a the pilot of a space shuttle climate data coming from different sources in a unified way in a single data screen through the *Pilot Climate Data System* (PCDS).
- Based on two fundamental ideas:
  - Including into the data file a **metadata** section, describing what the file contains.
  - An API between the application and the OS which decouples the data semantics and their storage logic.



# Scientific data libraries

## NetCDF

- NetCDF appears in 1990 as a spin-off of the CDF project.
- It maintains the fundamental CDF goals, plus:
  - Portability and platform-independent format.
  - C interface (CDF supported Fortran only).
  - UNIX implementation.
  - Mechanism to access multiple data points in a single step (hyperslab, same as using *slice* notation in Python).
- Common API used by every client: decouples data producers and consumers.



# Scientific data libraries

## NetCDF

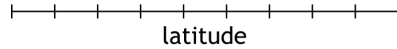
- A NetCDF file models a dataset using the following concepts:
  - Multidimensional variables:
    - Data types.
    - Shape (dimensionality).
    - Properties (key/value pairs).
  - Dimensions (coordinates):
    - Size.
    - Properties (key/value pairs).
  - File properties (key/value pairs).



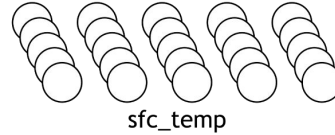
# Scientific data libraries

## NetCDF

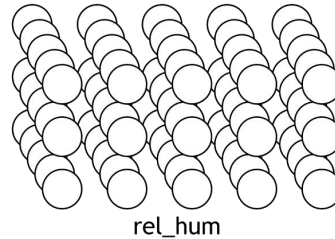
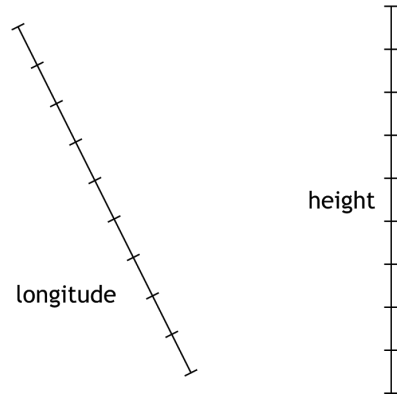
### Dimensions



### Variables



```
sfc_temp properties
long_name: "surface temperature"
units: "degC"
valid_range: -100, 70
```



```
rel_hum properties
long_name: "relative humidity"
units: "percent"
missing_value: -99
```



# Scientific data libraries

## NetCDF

- Along with the data, the NetCDF file includes:
  - The names of dimensions, variables, and properties.
  - The coordinate systems associated to variables.
  - The magnitudes of coordinates and variables.
  - The properties which capture any other descriptive information about the data.
- High-level interface: data are accessed using their names and the slice to access (instead of strings of bytes).
- Supports the basic scientific data types: strings and numbers, both integer and floating point.



# Scientific data libraries

## NetCDF

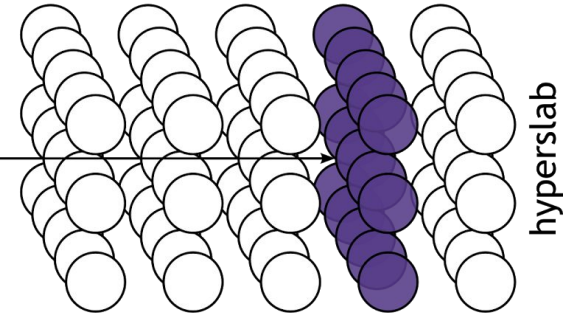
file id.

variable id.

corner  
vector

edge length  
vector

```
nc_get_vara_<type>( netcdf_id, var_id, start, count )  
nc_set_vara_<type>( netcdf_id, var_id, start, count )
```



<type> : text, uchar, schar, short, int, long,  
float, double, ushort, uint, longlong,  
ulonglong, string, (empty)

# Scientific data libraries

## NetCDF

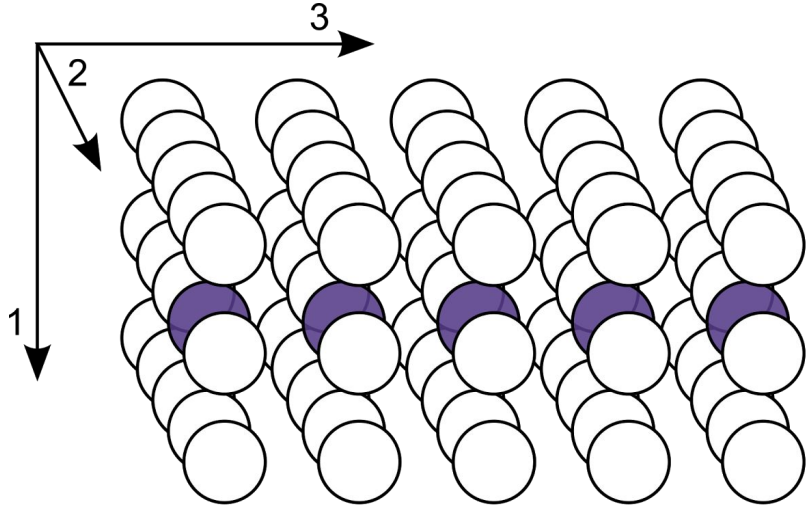
- We will use the Python interface.
- Instead of C-style functions (`nc_get_var()`, `nc_put_var()`) NumPy-style syntax is used, together with functions for creating and reading files.
- The NumPy-style slice syntax is readily translatable to hyperslab syntax:

$$[ \dots , e_k + a_k , \dots ] = [ \dots , e_k : e_k + a_k , \dots ]$$

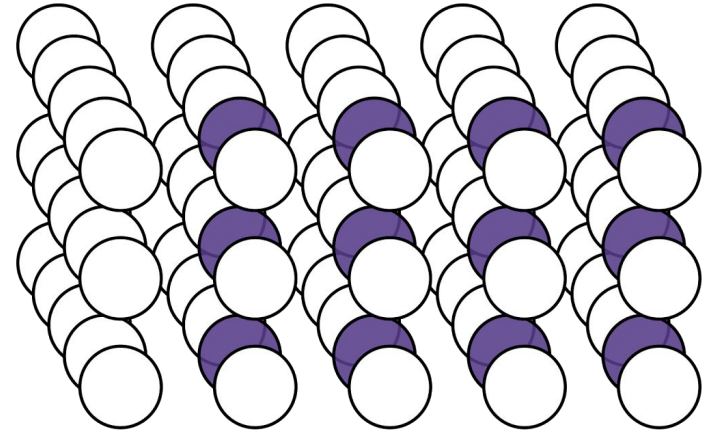


# Scientific data libraries

## NetCDF



corner : (1, 3, 0)  
edge length: (1, 1, 5)  
slices : (1, 3, :5)

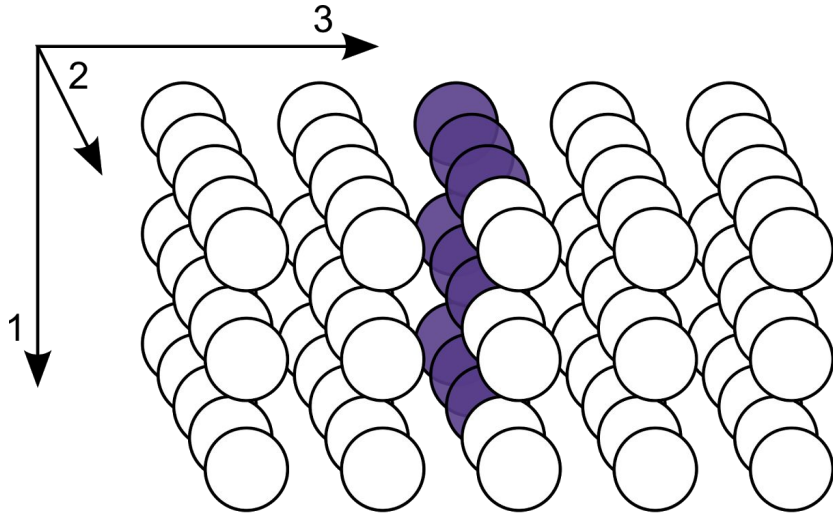


corner : (0, 3, 1)  
edge length: (3, 1, 4)  
slices : (:3, 3, 1:5)

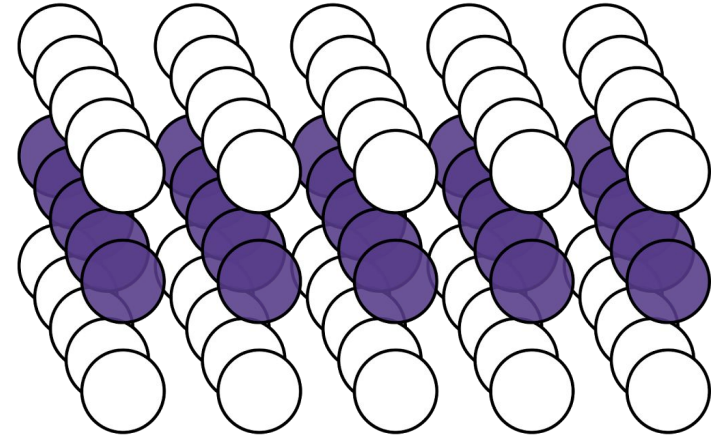


# Scientific data libraries

## NetCDF



corner : (0, 0, 2)  
edge length: (3, 3, 1)  
slices : (:3, :3, 2)

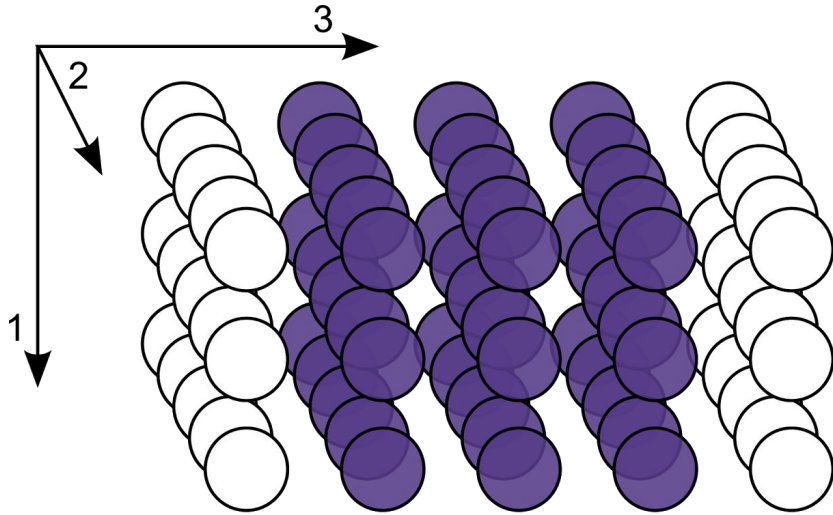


corner : (1, 0, 0)  
edge length: (1, 5, 5)  
slices : (1, :, :)

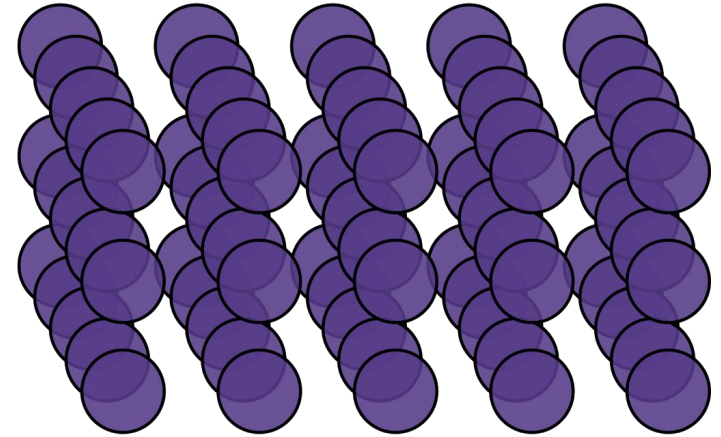


# Scientific data libraries

## NetCDF



corner : (0, 0, 1)  
edge length: (3, 5, 3)  
slices : (:, :, 1:4)



corner : (0, 0, 0)  
edge length: (3, 5, 5)  
slices : (:, :, :)



# Scientific data libraries

## NetCDF

- Hyperslab accesses in NetCDF always return contiguous memory, even if the original data were not contiguous on disk.
- It is not possible to use a “step” different from 1 in hyperslab access (less versatile than slices).
- NetCDF APIs are consistent with the memory conventions of the host language: e.g., the C interface is row-major, while the Fortran interface is column-major.
- A Fortran array defined as  $[N] [M]$  will be read in C/Python as  $[M] [N]$ .
- A single dimension can be unlimited (i.e., allowed to grow indefinitely after being created). It must be the slowest-changing dimension (outermost).



# Scientific data libraries

## NetCDF

- Internally, NetCDF (except version 4) uses the XDR format (*eXternal Data Representation*):
  - Non-proprietary format by Sun Microsystems to represent data in a platform-independent way.
  - Although XDR supports representing arbitrary C structs, NetCDF only supports multidimensional arrays of basic data types.
  - All conversions between machine-dependent formats are internally handled by XDR, in a programmer-transparent way. The application always receives data in the format native to the execution platform.



# Scientific data libraries

## NetCDF

Notebook-011.ipynb

- NetCDF-3 limitations:
  - Single unlimited dimension, due to internal representation limitations.
  - Not possible to perform sequential accesses to data: impossible to read/write from/to standard input/output. Intermediate files must be used.
  - Not a DB system. Impossible to access data except through names and coordinates (e.g., not possible to search all points of data with value greater than a given threshold).
  - Not possible to represent anything but multidimensional arrays (e.g., lists of structs).



# Scientific data libraries

## COARDS conventions

Property	Description
<code>long_name</code>	Descriptive name of a variable.
<code>scale_factor</code>	If present, the program reading data must multiply it by a scaling factor upon read.
<code>add_offset</code>	If present, data must be shifted by a specific amount upon read.



# Scientific data libraries

## COARDS conventions

Property	Description
<code>valid_range</code>	Valid range of values.
<code>_FillValue</code>	If present, the value will be used to initialize data on disk. It will be considered as a mark of NA values, so it should be outside of the valid range of values.
<code>missing_value</code>	Same as <code>_FillValue</code> , but with no special treatment (i.e., data are not automatically initialized on disk).



# Scientific data libraries

## HDF5

- HDF: Hierarchical Data Format.
- Created in 1987 in parallel to NetCDF.
- Developed by the University of Illinois at Urbana-Champaign and the National Center for Supercomputing Applications.
- Different objectives from NetCDF:
  - Platform-independent representation.
  - Access and store large objects efficiently.
  - Store different types of objects into a single container.
  - **Define a flexible format capable of accomodating new types of objects and metadata.**
  - Provide C and Fortran interfaces.



# Scientific data libraries

## HDF5

- During the 90s, and through NSF funding, HDF improves its end-user support: documentation, support, etc. and introduces compatibility with NetCDF files.
- In 1996 HDF5 is created on top of HDF4, providing parallel I/O for use in the Sandia, Los Alamos, and LLNL supercomputers.
- An HDF5 file is conceptually similar to a NetCDF file: it holds data and metadata, but working only with two types of objects:
  - Group: structure which aggregates zero or more HDF5 objects, plus metadata.
  - Dataset: Multidimensional data array, plus metadata.



# Scientific data libraries

## HDF5

- The hierarchy formed by groups and datasets is similar to a directory tree.
- Objects are referenced through paths:
  - `"/` is the root group in the file.
  - `"/foo` is object *foo* inside the root group in the file.
  - `"/foo/bar` is object *bar* inside group *foo* inside the root group in the file.
- An HDF5 dataset aggregates the concepts of variable and dimension in NetCDF: dimensions do not exist independently from a dataset.



# Scientific data libraries

## HDF5

- A dataset is created from:
  - Its path: a group to contain the object.
  - A name.
  - A datatype.
  - A *dataspace*, which encodes dimensionality.
- The `h5py` module allows to store in HDF5 any Python object, except for Unicode strings and generic objects.

# Scientific data libraries

## HDF5

- A dataspace defines the dimensionality of a dataset:
  - Two types: simple and complex:
    - Simple dataspace: multidimensional array.
    - Complex dataspace: different, more general organization.
  - Each dimension can have fixed or unlimited size. HDF5 allows multiple unlimited dimensions.
- When an HDF5 is created, its properties can be specified. Properties may refer to:
  - The HDF5 file.
  - Each dataset.
  - Data access.



# Scientific data libraries

## HDF5

- File creation properties:
  - User block size: a block at the beginning of the file which has free format, and can be written by the user at will (by default 0 bytes).
  - Offset and length pointer sizes (used in indices, by default `sizeof(size_t) = 8 bytes`).
  - Size of pointers in symbol table (by default 16 bits).
  - Size of pointers in B-trees to index chunked datasets (by default 32 bits).



# Scientific data libraries

## HDF5

- Dataset creation properties:
  - Storage type: contiguous, compact, chunked.
  - Data filters: compression, checksums, user-defined.
  - Memory allocation: early, incremental, late.
  - Fill values: ifset, alloc, never.
- Data access properties:
  - Caches and buffer sizes.
  - Metadata block sizes.
  - Parallel I/O mode.
  - etc.



# Scientific data libraries

## HDF5

Notebook-011.ipynb

- HDF5 associates an *attribute* list to groups and datasets.
- The attribute is implemented as a small dataset with its own dataspace.
- Differences with a regular dataset:
  - Does not accept partial I/O (hyperslabs).
  - Its datatype must match the memory type.



# Scientific data libraries

## Integration of NetCDF and HDF5

- NetCDF-4 is a joint project by Unidata and the HDF Group.
- It integrates both formats, using HDF-5 as the persistence layer of a NetCDF interface:
  - Allows to represent hierarchical data.
  - Preserves the simplicity of the NetCDF interface (this is only noticeable using the C and Fortran interfaces, not the Python one).
- NetCDF-4 is still compatible with previous NetCDF versions.
- It is capable of creating complex dataspace, as well as groups, multiple unlimited dimensions, compression, and parallel I/O.



# PyTables

- PyTables is a hierarchical database engine based in HDF5.
- Designed to efficiently managed large volumes of data.
- Optimizes memory and disk usage so that final size is smaller (particularly if compression is enabled) than using relational or object databases.
- Allows to define **tables** similar to those in a relational database.
- But each column can be a scalar or a NumPy array (i.e., multidimensional).
- A PyTables file stores both tables and NumPy arrays.



# PyTables

## Characteristics

- Supports up to  $2^{63}$  rows per table.
- Unlike relational databases, each cell in a hierarchical table can have any dimensionality.
- Supports indexing any column of the table.
- Users can define metadata that give semantical meaning to their data.
- Supports data compression through Zlib, LZO, bzip2, and Blosc.
- Inherits HDF5 characteristics: supports files larger than 2GB, high-performance I/O, platform-independent files, ...



# PyTables

## Characteristics

Notebook-012.ipynb

- All the objects of a PyTables tree are instances of the `Node` class.
- The `Group` and `Leaf` classes inherit from `Node`: a group contains other groups and leaf objects (similar to HDF5).
- A leaf object is a data container, and can be an instance of the following classes:
  - `Table`.
  - `Array`.
  - `CArray`: *chunked*, allows data compression.
  - `EArray`: *extensible*, allows to enlarge one of its dimensions.
  - `VLAArray`: *variable length*, its rows have variable number of columns.
  - `UnImplemented`: HDF5 data unsupported by PyTables.



# Blaze

- NumPy arrays and Pandas `DataFrames` are the foundational structures of Python as a data analysis tool.
- However, processing capabilities are restricted by physical memory.
- Besides, loading data into Pandas is not always trivial (e.g., parsing JSON, XML, SQL, etc.).
- Blaze has two main goals:
  - To integrate different computational backends and data storage technologies transparently.
  - To trivialize out-of-core computing.



# Blaze

## Data sources

- Native Python objects (lists, dictionaries, ...).
- NumPy.
- Pandas.
- SQL.
- HDF5.
- MongoDB.
- Spark.
- Impala.



# Blaze

- Blaze does not compute anything.
- It “simply” performs data transformations and bridges software interfaces to make different data sources and computational backends interact in a cohesive way.
- It employs functionalities which already exist in the Python ecosystem.
- Allows to conveniently transform data from one format to another (e.g., CSV to SQL).
- Blaze is not as mature as Pandas. Its documentation is often incomplete, and its APIs change frequently.



# Blaze

Notebook-013.ipynb

- A Blaze *server* provides a web-services API to data.
- Any data source accepted by Blaze can be exposed through this API.
- The server can be accessed in different ways:
  - Using programs written in any language, e.g., using the `requests` module, to build JSON calls.
  - Using the `blaze` module, defining a data source with the server address through the `blaze://` protocol.



# Seaborn

Notebook-014.ipynb

- Seaborn is a graphics library which 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.



# Bokeh

Notebook-015.ipynb

- Bokeh is another graphical library with different objectives.
- Part of Continuum's (Anaconda) analysis tools, together with Blaze.
- Focuses on building plots that can be visualized as HTML with Javascript.
- Creates *interactive* plots.
- Provides output in different formats, such as HTML, inlined into an iPython Notebook, or PNG.
- Does not currently allow to create vector versions (e.g., SVG).



# Performance analysis

- Python is intrinsically inefficient when executing array codes (i.e., scientific/engineering codes).
- (Except if said codes can be “vectorized” using NumPy).
- However, it is very good at I/O-bound codes such as data analytics.
- Ideally, we could use Python as an universal language, including:
  - Data generation (main computation).
  - Data analytics (statistics, *crunching*).
  - Report generation (plots).



# Performance analysis

- We focus now on implementing efficient scientific codes on Python.
- The development cycle for scientific codes is generally incremental: analysis, implementation, debugging, evaluation, rinse and repeat.
- During performance analysis it is fundamental to understand the bottlenecks and their causes to correct them in a future iteration.
- Profiling is the process which investigates the behavior of the code, developing a model which allows to correct the observed problems in successive development phases.



# Performance analysis

## Measuring runtime

Notebook-016.ipynb

- The execution time of a section of code is a useful performance metric.
- The `time` module allows to measure the execution time of a block of code.
- iPython also provides the `%time` and `%timeit` magic commands.
- However, we need to investigate how execution time is distributed across the code.
- It is not convenient to include a great deal of profiling code including calls to `time` and/or `print`.
- Python provides powerful tools to analyze code performance in an automatic way, avoiding (most) changes to the code.



# Performance analysis

## cProfile

Notebook-016.ipynb

- The `cProfile` module, part of the Python standard library, interacts with the Python VM to measure the time spent inside each function in a code.
- Its output is basically identical to that of GNU `prof`: a list of functions including number of times each one was called, cumulative execution time, and time per call.
- Due to the required instrumentation, executing code through `cProfile` is much slower than a regular execution.
- Information is provided **at the function level** only, i.e., no information can be inferred about which parts of a function are more costly.



# Performance analysis

## `line_profiler`

Notebook-016.ipynb

- The `line_profiler` module analyzes the time spent executing each line of a function.
- Its output is similar to that of `cProfile`, but instead of offering data at the function level, it does so at the line-of-code level.
- As any other profiling tool, it introduces execution overheads. In this case, much larger to those of `cProfile`.
- For this reason, the ideal chain of tools is: analyze the full code with `cProfile`, and then do a fine-grained analysis of the critical functions in the code with `line_profiler`.



# Compiling to C: Cython

Notebook-016.ipynb

- Cython is a framework allowing to compile Python code into C.
- A Cython code is similar to OpenCL: there is a main program written in Python which calls computational kernels written in annotated Python (`.pyx` files).
- This `.pyx` file is translated to C by a source-to-source compiler.
- The C code obtained from an un-annotated `.pyx` is very similar to the original Python code: it calls the Python interpreter to gain dynamic knowledge about data types and function calls.
- If the `.pyx` with the types of the variables in the code it is possible to save runtime checks, and therefore interpreter calls, getting dramatic performance returns with minimal effort.



# Compiling to C

## Parallelism in Python

Notebook-016.ipynb

- One of the most controversial characteristics of Python is the *Global Interpreter Lock* (GIL), which prevents two different threads to execute Python bytecode at the same time.
- The GIL has several advantages (single-threaded execution is more efficient, integration with C libraries, etc.), but it clearly does not benefit parallelism.
- **It does not affect I/O-bound operations.** A process releases the GIL before waiting for I/O.
- Cython allows automatic parallelization through OpenMP.
- The limitations imposed by the GIL are worked-around by working with native C types, avoiding calls to the interpreter.



# Compiling to C

## Numba

Notebook-017.ipynb

- Numba is a tool similar to Cython: it generates C code from Python code using *Just-In-Time compilation* (JIT).
- Besides their interfaces, the main difference between Numba and Cython is the framework used for source-to-source compiling. It uses its own LLVM implementation.
- Numba provides the `@jit` decorator, as well as `numba.jit()`, to compile particular functions.
- The paid NumbaPro version also allows to automatically generate GPGPU code.



# Parallel Python

## threading

Notebook-018.ipynb

- The `threading` module allows to write multithreaded programs.
- It works well for I/O-bound codes, but not in CPU-bound ones except when the computation uses C libraries intensively (such as Cython or NumPy).
- Parallelism model similar to POSIX threads (`pthread`s).



# Parallel Python

## multiprocessing

Notebook-018.ipynb

- The `multiprocessing` module imitates the `threading` interface, but creates processes instead of threads.
- Each process executes on its own Python interpreter, therefore circumventing the single-GIL limitation.
- However, different processes do not share memory. The `Manager` class allows to share different types of Python objects (lists, dictionaries, mutexes, queues, arrays, ...).

