

# Python for Finance

## Lecture Notes

Matteo Sani

Quants Staff - MPS Capital Services  
[matteo.sani@mpscapitalservices.it](mailto:matteo.sani@mpscapitalservices.it)



# Contents

<b>1</b>	<b>Introduction to python</b>	<b>5</b>
1.1	What is python ? . . . . .	5
1.2	Python basics . . . . .	8
1.2.1	Variables . . . . .	8
1.2.2	Boolean expressions . . . . .	10
1.2.3	String expressions . . . . .	11
1.2.4	Mathematical expressions . . . . .	12
1.3	Modules . . . . .	14
1.4	Indented blocks and the if/elif/else statement . . . . .	17
1.5	Loops . . . . .	18
1.5.1	for . . . . .	18
1.5.2	while . . . . .	20
<b>2</b>	<b>Data Containers</b>	<b>21</b>
2.1	Lists . . . . .	21
2.2	Dictionaries . . . . .	24
2.3	Tuples . . . . .	27
<b>3</b>	<b>Date and Time</b>	<b>29</b>
3.1	Dates . . . . .	29
<b>4</b>	<b>Python's Object Oriented Programming</b>	<b>31</b>
4.1	Functions . . . . .	31
4.2	Variable scope . . . . .	33
4.3	Classes . . . . .	36
4.3.1	Inheritance and Overriding Methods . . . . .	38
<b>5</b>	<b>Data Manipulation and Its Representation</b>	<b>41</b>
5.1	Getting Data . . . . .	41
5.1.1	Manage Data . . . . .	43
5.1.2	Unwanted observations and outliers . . . . .	43
5.1.3	Handle Missing Data . . . . .	46
5.1.4	Filter, Sort and Clean Data . . . . .	46
5.2	Plotting in python . . . . .	50
5.2.1	Plot a graph given $x$ and $y$ values (scatter-plot) . . . . .	50
5.2.2	Plotting an Histogram . . . . .	52

<b>6</b>	<b>Interpolation, Discount Factors and Forward Rates</b>	<b>55</b>
6.1	Linear interpolation . . . . .	55
6.1.1	Log-linear interpolation . . . . .	57
6.1.2	Limitations of Interpolation . . . . .	59
6.2	Discount curve interpolation . . . . .	59
6.3	Forward Rates . . . . .	64
6.3.1	2008 Financial Crisis . . . . .	65
<b>7</b>	<b>Swaps and Bootstrapping</b>	<b>69</b>
7.1	Payment Dates Generator . . . . .	69
7.2	Overnight Index Swap . . . . .	69
7.2.1	OIS Valuation . . . . .	70
7.2.2	OvernightIndexSwap Class . . . . .	71
7.3	Bootstrap Technique . . . . .	73
7.3.1	Building OIS instances . . . . .	73
7.3.2	Constructing the Yield Curve . . . . .	75
7.3.3	Bootstrap as Minimization Problem . . . . .	77
7.3.4	Minimization Algorithm . . . . .	77
7.3.5	OIS Example . . . . .	80

# Chapter 1

## Introduction to python

Python is one of the most widely used programming languages in the world, and it has been around for more than 28 years now.

First and foremost reason why python is much popular because it is highly productive as compared to other programming languages like C++ and java. It is a much more concise and expressive language and requires less time, effort, and lines of code to perform the same operations.

This makes python very easy-to-learn programming language even for beginners and newbies. It is also very famous for its simple programming syntax, code readability and English-like commands that make coding in python lot easier and efficient. With python, the code looks very close to how humans think. For this purpose, it must abstract the details of the computer from you. Hence, it is slower than other “lower-level language” like C.

There were times when computer run time was to be the main issue and the most expensive resource. But now, things have changed. Computer, servers and other hardware have become much much cheaper than ever and speed has become a less important factor. Today, development time matters more in most cases rather than execution speed. Reducing the time needed for each project saves companies tons of money.

As far as the execution speed or performance of the program is concerned, we can easily manage it by horizontal scaling, meaning that more servers can be used to reach that level of speed or performance.

In short, python is widely used even when it is somehow slower than other languages because:

- is more productive;
- companies can optimize their most expensive resource: employees;
- rich set of libraries and frameworks;
- large community.

### 1.1 What is python ?

---

Python is a so called *interpreted language*: it takes some code (a sequence of instructions), reads and executes it. This is different from other programming languages like C or C++ which *compile* code into a language that the computer can understand directly (*machine language*).

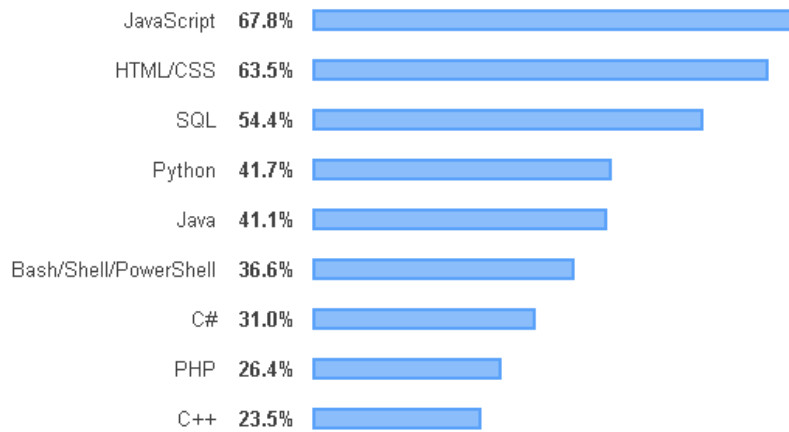


Figure 1.1: List of most used languages by developers according to Stack Overflow survey in 2019.



Figure 1.2: List of most loved languages by developers according to Stack Overflow survey in 2019.



Figure 1.3: List of most *dreaded* languages by developers according to Stack Overflow survey in 2019.

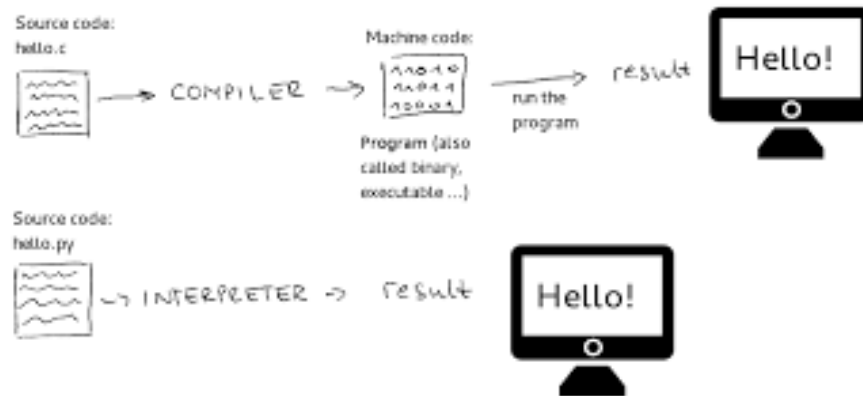


Figure 1.4: Interpreted vs compiled language

As a result, python is essentially an *interactive* programming language, you can program and see the results almost at the same time. This is very nice for a faster development since compilation time can be quite long (just to give an idea the compilation of our C++ financial code takes more than one hour). However there are drawbacks in terms of performance, the *translation* to machine language has to be done in real-time resulting in slower execution times.

High-level program	<pre> class Triangle {     ...     float surface()         return b*h/2; } </pre>
Low-level program	<pre> LOAD r1,b LOAD r2,h MUL r1,r2 DIV r1,#2 RET </pre>
Executable Machine code	<pre> 0001001001000101 0010010011101100 10101101001... </pre>

Figure 1.5: Human readable vs machine code

In the next chapters we'll take a quick tour of python and see the main features and characteristics of this programming language, later on we will see how it can be useful to solve real-world financial problems.

First of all since python, as basically all programs, comes in different version and flavours we need to specify the particular one we are going to use. The latest version (at the time I'm writing this pages) is 3.8.5, but it is continuously evolving, however it is not difficult to see older versions floating around (e.g. 2.7). This is because there are some big differences between python2.X and python3.X which prevent a sizable portion of python2 users to stick with it (consider that moving to python3 would require a large amount of work to adapt big projects). In conclusion we will concentrate on python 3.7.

## 1.2 Python basics

Every language has *keywords*, these are reserved words that have a special meaning and tell the computer what to do. The first one we encounter is `print`: it prints to screen whatever is specified between the parenthesis.

```
print ("Hello world !")
```

```
Hello world !
```

```
print ("Welcome")  
print ("to")  
print ("everybody")
```

```
Welcome  
to  
everybody
```

Good programming practice recommends to document the code you write (you will soon see that it is surprisingly easy to forget what you wanted to do in your code). In python you can add comments to code starting a line with a hash character (#).

```
print ("Ciao") # this is a comment
```

```
Ciao
```

### 1.2.1 Variables

A variable is a computer memory location paired with a symbolic name, which contains some quantity of information referred to as a *value* (e.g. a number, a string...). Variables and hence data they contain, can be used, referenced and manipulated throughout a program. A value is assigned to a variable with the equal operator (=) and printing a variable shows its content.

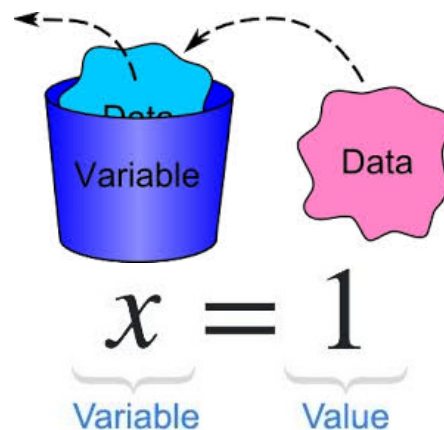


Figure 1.6: Graphical representation of a variable.



```
x = 9  
print (x)
```

9

```
myphone = "Huawei P10Lite"  
print (myphone)
```

Huawei P10Lite

Another very useful keyword is `type`: it tells which kind of object is stored in a variable.

```
print (type(x))
print (type(myphone))

<class 'int'>
<class 'str'>
```

After their definitions `x` and `myphone` can be used as aliases for a number and a string and their content manipulated, for example:

```
print (x+5)

14
```

There are rules that limit the variable naming possibilities, in particular they must:

- begin with a letter (`myphone`) or underscore (`_myphone`);
- other characters can be letters, numbers or more `_`;
- variable names are case-sensitive so `myphone` and `myPhone` are two distinct variables;

**Keywords, as said, are reserved words and as such cannot be used as variable names (e.g. `print`, `type`, `for`...).**

To use **good** variable names (and make your programs clearer and easier to read) always choose meaningful names instead of short names (i.e. `numberOfCakes` is much better than simply `n`), try to be consistent with your conventions (e.g. choose once and for all between `number_of_cakes` or `numberofcakes` or `numberOfCakes`, usually begin a variable name with underscore (`_`) only for a special case (will see later when this is usually done).

### 1.2.2 Boolean expressions

Boolean expressions evaluate to `true` or `false` only. This type of expressions usually involve logical or comparison operators like `or`, `and`, `>` (greater-than), `<` (less-than),... The equal-to Boolean operator symbol is a double `=` (`==`), to not be confused with the assignment operator single `=` (`=`), with the first we compare two variables, with the second we associate a value to a variable.

Let's see some example. The following expression answers the question is 1 equal to 2:

```
1 == 2

False
```

Here another example using the not equal operator (`!=`):

```
1 != 2

True
```

```
2 < 2

False
```

```
2 <= 2 # in this case we allow the numbers to be equal too
```

```
True
```

```
print (x)
```

```
15 <= x and x <= 20 # this expression could also be written as 15 <= x <= 20
```

```
11
```

```
False
```

```
15 <= x or x <= 20
```

```
True
```

```
not (x > 20) # the not keyword negates the following expression
```

```
True
```

### 1.2.3 String expressions

A “string” is a sequence of characters (letters, digits, spaces, punctuation,...). There are many operations that can be performed on strings, like for example concatenate (with + operator), truncate, replace characters,...

```
mystring = "some text with punctuation, spaces and digits 10"
```

```
mystring.replace("s", "z")
```

```
'zome text with punctuation, zpacez and digitz 10'
```

```
"abc" + "def" # it is possible to concatenate strings with +
```

```
'abcdef'
```

```
"The number " + 4 + " is my favourite number"
```

```
# this causes an error since we are trying to concatenate a string
```

```
# with a number so two different kind of objects
```

```
-----
```

```
TypeError
```

```
Traceback (most recent call last)
```

```
<ipython-input-33-b9f65c5a45f7> in <module>()
```

```
----> 1 "The number " + 4 + " is my favourite number"
```

```
      2 # this causes an error since we are trying to concatenate a string
```

```
      3 # with a number so two different kind of objects
```

```
TypeError: can only concatenate str (not "int") to str
```

To avoid this error is possible to **cast** an object to a different type which means to convert an object to a different type. In this case we can *force* the number four to be represented as a string with the `str()` function:

```
"The number " + str(4) + " is my favourite number"

'The number 4 is my favourite number'
```

```
print (type(3.4))
print (type(str(3.4)))

<class 'float'>
<class 'str'>
```

In this simple case everything worked fine but type casting is not always possible: for example a number can be converted to a string (e.g. from the integer 4 to the actual symbol “4”) but the opposite is not possible (e.g. cannot convert the string “matteo” to a meaningful number). In this second case we can try to use the function `int()` to convert a string to an integer.

```
int("matteo")

-----

ValueError                                Traceback (most recent call last)

<ipython-input-17-979283bb65e4> in <module>
----> 1 int("matteo")

ValueError: invalid literal for int() with base 10: 'matteo'
```

```
int("4")

4
```

### Pretty string formatting

In order to get prettier strings than those obtained just concatenating with the `+` operator, python allows to format text using the following syntax `"text {} other text {}".format(var1, var2)`. With this notation, each `{}` is mapped to the variables listed in the format statement, the optional characters inside the curly brackets can determine the resulting format, for example in the following code `{:.1f}` means that this variable is a float number and that has to be printed with only one digit only after the decimal separator.

```
"The speed of light is about {:.1f} {}".format(299792.458, "km/s")

'The speed of light is about 299792.5 km/s'
```

In addition format allows for 0-padding of numbers, left or right alignment of text columns and so on.

### 1.2.4 Mathematical expressions

Below few examples of the basic mathematical expressions available in python.

```
1 + 2
```

```
3
```

```
40 - 5
```

```
35
```

```
x * 20 # remember that we set x equal to 9
```

```
180
```

```
x / 4
```

```
2.25
```

```
print (type(2.25))
```

```
<class 'float'>
```

```
x // 4 # integer division - result will be truncated to the  
      # corresponding integer (no rounding)  
      # 11 / 3 = 3.666666 -> 11 // 3 = 3
```

```
2
```

```
y = 3
```

```
x ** y # x to the power of y
```

```
729
```

```
3 * (x + y)
```

```
36
```

As an example of variable manipulation let's try to increment x by 1 and save the result again in x.

```
print (x)
```

```
x = x + 1
```

```
print (x)
```

```
15
```

```
16
```

Sometimes the increment of a variable plus the assignment to the same variable is written with a more compact syntax `x += 1` (this is also true for other operators e.g. `x *= 2`).

More complex mathematical functions are not directly available, let's see for example the logarithm:

```
log(3)
```

```
-----
```

```
NameError                                Traceback (most recent call last)

<ipython-input-17-ffde4d60496a> in <module>()
----> 1 log(3) # causes an error because the logarithm function
      2      # is not available by default

NameError: name 'log' is not defined
```

## 1.3 Modules

One very important feature of each language is the ability to reuse code among different programs, e.g. imagine how awful would be if you had to re-implement every time you need it a function to compute the logarithm. Usually there are mechanisms that allow to collect useful routines in *packages* (or *libraries*, or *modules*) so that later they can be called and used by any program may need them.

These collections of utilities in python are called *modules* and each installation of this language brings with it a standard set of them. If you need more functionality, you can download more modules from the web (there are zillions out there) or if you are not satisfied with what you found you can write your own (which is one the goal of this course in the end).

Some examples of useful modules we will use are:

- Numpy - which provides matrix algebra functionality and much more;
- Scipy - which provides a whole series of scientific computing functions;
- Pandas - which provides tools for manipulating time series or data-set in general;
- Matplotlib - for plotting graphs;
- Jupyter - for notebooks like this one.

Later we will take a closer look at three modules which are quite useful in financial analysis.

In order to load a module in a python program you can use the `import` keyword. To inspect a module (to understand which are its functionalities) it can be used the `help` and `dir` keywords: the first write a help message which usually describes the functionalities of a module, the latter list all the available functions of a module. **In order to access a function of a module you have to use the . (dot) operator:** `module-name.function-name`.

Let's see an example dealing with the `math` module which implements the most common mathematical functions.

```
import math
dir(math)

Out[18]: ['__doc__',
          '__loader__',
          '__name__',
          '__package__',
          '__spec__',
          'acos',
          'acosh',
```

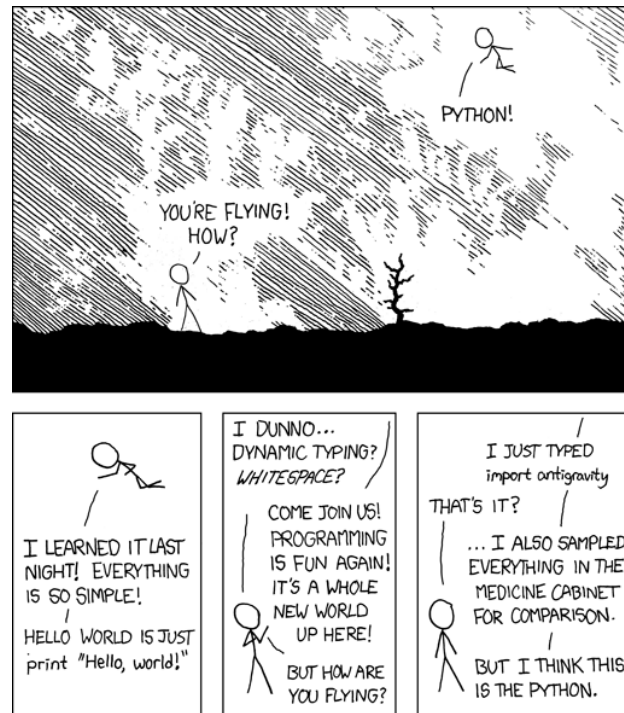


Figure 1.7: Python has many modules for download on the web...

```
'asin',
'asinh',
'atan',
'atan2',
'atanh',
'ceil',
'copysign',
'cos',
'cosh',
```

```
...
```

```
help(math)
```

```
Help on module math:
```

```
NAME
```

```
    math
```

```
MODULE REFERENCE
```

```
    https://docs.python.org/3.6/library/math
```

The following documentation is automatically generated from the Python source files. It may be incomplete, incorrect or include features that are considered implementation detail and may vary between Python implementations. When in doubt, consult the module reference at the

location listed above.

#### DESCRIPTION

This module is always available. It provides access to the mathematical functions defined by the C standard.

#### FUNCTIONS

```
acos(...)
    acos(x)
```

Return the arc cosine (measured in radians) of x.

...

```
math.log(3)
```

```
1.0986122886681098
```

```
math.exp(3)
```

```
20.085536923187668
```

```
print (type(math.log)) # yet another type: builtin function
print (type(math.log(3)))
```

```
<class 'builtin_function_or_method'>
<class 'float'>
```

If we want to avoid to type `math.` every time we compute a logarithm or an exponential, we can just import only the needed functions from a module using the following syntax:

```
from math import log, exp
print (log(3))
print (exp(3))
```

```
1.0986122886681098
```

```
20.085536923187668
```

As an example let's compute the interest rate  $r$  that produces a return  $R$  of 11000 Euro when investing 10000 Euro for 2 years:

$$R = Ne^{r\tau} \rightarrow r = \frac{1}{\tau} \log\left(\frac{R}{N}\right)$$

```
rate = (1/2)*log(11000/10000)
print (rate)
```

```
0.04765508990216247
```



## 1.4 Indented blocks and the if/elif/else statement

Unlike other languages which uses parenthesis to isolate blocks of code python uses *indentation*. A first example of this peculiarity is given by if/elif/else statements. Such commands allow to dynamically run different blocks of code based on certain conditions. For example in the following we print different statements according to the value of x, note that the block of code to be run according each condition is shifted (i.e. indented) with respect to the rest of the code:

```
print (x)
if x == 1:
    print ("This will not be printed")
    # the block of code that is run if the first condition is met is indented
elif x == 15:
    print ("This will not be printed either")
    # again the block of code that is run here is indented
    # to be "isolated" by the rest
else:
    print ("This *will* be printed")

16
This *will* be printed
```

If by mistake the indentation of a block is missing an error is raised:

```
if x == 1:
    print ("This will not be printed")
elif x == 15:
    print ("This will not be printed either")
else:
    print ("This *will* be printed")

File "<ipython-input-38-4535a45a6419>", line 3
    print ("This will not be printed")
    ^
IndentationError: expected an indented block
```

Below another example:

```
if x != 1:
    print ("x does not equal to 1")

x does not equal to 1
```

Just for comparison this is the same code written in C++:

```
if (x == 1) {
    print ("This will not be printed");
}
else if (x == 15) {
    print ("This will not be printed either");
}
```

```
}  
else {  
print ("This *will* be printed");  
}
```

N.B. Notice how indentation doesn't matter at all here since the blocks are enclosed and defined by the brackets.

## 1.5 Loops

---

Another very important feature of a language is the ability to repeatedly run the same block of code many times. This is called looping and in python can be done with `for` or `while` keywords.

### 1.5.1 `for`

In a `for` loop we specify the set (or interval) over which we want to loop and a variable will assume all the values in that set (or interval). For example let's assume we want to print all the numbers between 25 and 30 excluded (here the keyword `range` returns the list of integers between the specified limits, if the first limit is not specified 0 is assumed):

```
for i in range(25, 30):  
    print (i)  
  
25  
26  
27  
28  
29
```

At each cycle of the loop the variable `i` takes one of the values between 25 and 31. With `range` it is also possible to specify the step, so that the loop can jump every 2 units or to go in descending order:

```
for i in range (30, 25, -1):  
    print (i)  
  
30  
29  
28  
27  
26
```

If it is needed to skip values in the loop the `continue` keyword can be used; in the code below 5 is actually missing from the list in the printout since it has been skipped by the `continue`:

```
for i in range(10):  
    if i == 5:  
        continue  
    print (i)
```

```
0
1
2
3
4
6
7
8
9
```

Instead of using range it is possible to specify directly the set of looping values:

```
for i in (4, 6, 10, 20): # here we loop directly on a list of numbers
    print (i)

4
6
10
20
```

Finally looping on a string actually means to loop on each single character:

```
phrase = 'how to loop over a string'
for c in phrase:
    print (c)

h
o
w

t
o

l
o
o
p

o
v
e
r

a

s
t
r
i
n
```

```
g
```

### 1.5.2 while

In a for loop we go through all the elements of a list of objects, the while statement instead repeats the same block of code until a condition is met. The following block of code is run if `x` squared is less than 50, so we first set `x=1` and at each iteration we increment it by 1 until the condition is True (8 squared is 64 which is greater than 50):

```
x = 1
while x ** 2 < 50:
    print (x)
    x += 1

1
2
3
4
5
6
7
```

It is possible to exit prematurely from a while loop using the `break` keyword. In this case the while-condition is simply `True` so the code would run forever unless we set an exit strategy.

```
x = 1
while True:
    if (x ** 2 > 50):
        break
    print (x)
    x += 1

1
2
3
4
5
6
7
```

## Chapter 2

# Data Containers

In this chapter the container types available in python are reviewed.

### 2.1 Lists

---

A list in python is a container that is a *mutable*, ordered sequence of elements. Each element or value that is inside of a list is called an *item*. Each item can be accessed using square brackets notation (very important, list indexing is zero-based so the first element has index 0 actually). A list is considered mutable since you can add, remove or update the items in it. Ordered instead means that items are kept in the same order they have been added. Lists can be created by enclosing in square brackets the comma-separated list of the items or using the `list()` operator.

```
mylist = list([21, 32, 15])
mylist = [21, 32, 15]
print(mylist)
print (type(mylist))

[21, 32, 15]
<class 'list'>
```

```
print(mylist[0])

21
```

If you have a list of lists (i.e. a 2-dimensional list) you can use the square brackets multiple times to access the inner elements:

```
alist = [[1,2], [3,4], [5,6]]
print (alist[1][1]) # first [1] returns [3,4], second returns 4
```

The number of elements in a list is counted using the keyword `len()`:

```
print(len(mylist))

3
```

Looping on list items can be achieved in two ways: using directly the list or by index:

```
print ("Loop using the list itself:")
for i in mylist:
    print (i)

print ("Loop by index:")
for i in range(len(mylist)): # len() returns the number of items in a list
    print (mylist[i])

Loop using the list itself:
21
32
15

Loop by index:
21
32
15
```

With the enumerate function is actually possible to do both at the same time since it returns two values, the index of the item and its value, so in the example below, i will take the item index values while item the item value itself:

```
for i, item in enumerate(mylist):
    print (i, item)

0 21
1 74
2 85
3 15
4 188
```

Since a list is mutable we can dynamically change its items:

```
mylist[1] = 74 # we can change list items since it's *mutable*
print (mylist)

[21, 74, 15]
```

With append an item is added at the end, while with insert an item can be added in a specified position:

```
mylist.append(188) # append add an item at the end of the list
print (mylist)

[21, 74, 15, 188]

mylist.insert(2, 85) # insert an item in the desired position
                    # (2 in this example)
print (mylist)

[21, 74, 85, 15, 188]
```

To append multiple values at once to a list a loop can be used but python offers a single line way of doing it: `[i*2 for i in range(10)]`. This syntax is called *list comprehension*.

Accessing items outside the list range gives an error:

```
mylist[10] # error ! it doesn't exists, the list has only 3
           # elements, so the last is item 2

-----

IndexError                                Traceback (most recent call last)

<ipython-input-36-ed1e5e6c3e46> in <module>
----> 1 mylist[10] # error ! it doesn't exists, the list has only 3
      2           # elements, so the last is item 2

IndexError: list index out of range
```

Read carefully the error messages usually they are very explicative and can help a lot in *debugging* (i.e. finding mistakes) in your programs.

There are two more nice features of python indexing:

- negative indices are like positive ones except that they starts from the last element;
- *slicing* which allows to specify a range of indices to select more items at once (if the first or last limits are missing slicing will start from the first or end with last index respectively).

```
print ("negative index -1 returns the last element:", mylist[-1])
print ("slice [1:3] returns items 1st and 2nd:", mylist[0:3])
print ("slice [:2] returns items 0th and 1st:", mylist[:2])
print ("slice [2:] returns items between the 2nd and the last:", mylist[2:])
```

negative index -1 returns the last element: 188  
slice [1:3] returns items 1st and 2nd: [21, 74, 85]  
slice [:2] returns items 0th and 1st: [21, 74]  
slice [2:] returns items between the 2nd and the last: [85, 15, 188]

Needless to say that slicing with `[:]` returns the entire list.

It is worth mentioning that a list doesn't have to be populated with the same kind of objects (list indices are instead always integers).

```
mixedlist = [1, 2, "b", math.sqrt]
print (mixedlist)
```

```
[1, 2, 'b', <built-in function sqrt>]
```

```
print (mixedlist['k'])
```

---

TypeError

Traceback (most recent call last)

```
<ipython-input-72-aea4c7f9789e> in <module>()
----> 1 print (mixedlist['k'])
```

```
TypeError: list indices must be integers or slices, not str
```

A complete list of the commands available for a list can be shown with the `dir` statement:

```
dir(list)

[...
'append',
'clear',
'copy',
'count',
'extend',
'index',
'insert',
'pop',
'remove',
'reverse',
'sort']
```

Their meaning is pretty clear, so for example `sort` re-order the items according to a custom criteria or `index(item)` return the index of the specified item.

## 2.2 Dictionaries

As we have seen lists are ordered collections of elements and as such we can say that map integers (the index of each item) to values (any kind of python object). *Dictionaries* generalize such a concept being containers which map *keys* (**almost** any kind of python object) to values (any kind of python object).

In our previous section we had:

```
0 (0th item) → 21
1 (1st item) → 74
2 (2nd item) → 85
...
```

With a dictionary we can have something like this:

```
"apple"(key) → 4
"banana"(key) → 5
```

As we will see dictionaries are very flexible and will be very useful to represent complex data structures.

Dictionaries can be created by enclosing in curly brackets the comma-separated list of key-value pairs (key and value are separated by a `:`), or using the `dict()` operator. In lists we could



access items by index, here we do it by key still using the square brackets. Trying to access not existing keys results in error, but we can check if a key exists with the `in` operator. As before, if a dictionary contains other dictionaries or lists, the square brackets can be applied repeatedly to access the inner items.

```
adict = {"apple": 4, "banana": 5}
print (adict["apple"])

4
```

```
adict["pear"] # error !

-----

KeyError                                Traceback (most recent call last)

<ipython-input-41-9d051ebd10de> in <module>
----> 1 adict["pear"] # error ! this key doesn't exists

KeyError: 'pear'
```

```
"pear" in adict # indeed

False
```

The items can be dynamically created or updated with the assignment `=` operator, while again `len()` returns the number of items in a dictionary.

```
adict["banana"] = 2
adict["pear"] = 10
print (len(adict))
print (adict)

3
{'apple': 4, 'banana': 2, 'pear': 10}
```

Dictionaries can be made of more complicated types than simple string and integers:

```
adict[math.log] = math.exp
```

Also dictionaries can be created with the *comprehension* syntax: `{i:v for i, v in enumerate(["a", "b", "c"])}`.

Looping over dictionary items can be done by key, by value or by both: `.keys()` returns a list of keys, `.values()` returns a list of values and `.items()` a list of pairs key-value.

```
print ("All keys: ", adict.keys())
for key in adict.keys():
    print (key)

print ("All values: ", adict.values())
for value in adict.values():
    print (value)
```

```

print ("All key-value pairs: ", adict.items())
for key, value in adict.items():
    print (key, value)

All keys: dict_keys(['apple', 'banana', 'pear', <built-in function log>])
apple
banana
pear
<built-in function log>

All values: dict_values([4, 2, 10, <built-in function exp>])
4
2
10
<built-in function exp>

All key-value pairs: dict_items([('apple', 4), ('banana', 2), ('pear', 10),
(<built-in function log>, <built-in function exp>)])
apple 4
banana 2
pear 10
<built-in function log> <built-in function exp>

```

To merge two dictionaries the function `update()` can be used, while with `del` it is possible to remove a key-value pair.

```

del adict[math.log]
seconddict = {"watermelon": 0, "strawberry": 1}
adict.update(seconddict)
print (adict)

{'apple': 4, 'banana': 2, 'pear': 10, 'watermelon': 0, 'strawberry': 1}

```

Again the complete list of dictionary functions can be shown with `dir`:

```

dir(dict)

[...
'clear',
'copy',
'fromkeys',
'get',
'items',
'keys',
'pop',
'popitem',
'setdefault',
'update',
'values']

```

## 2.3 Tuples

Tuples create a bit of confusion for beginners because they are very similar to lists but they have some subtle conceptual differences. Nonetheless, tuples do appear when programming in python so it's important to know about them.

```
list1 = [1,2,3,4]
```

- List

```
tuple1 = (1,2,3,4)
```

- Tuple

Figure 2.1: At first glance list and tuples look very similar, but they are not...

Like lists, tuples are containers of any type of object. Unlike lists though they are *immutable* which means that once they have been created the content cannot be changed (i.e. no append, insert or delete of the elements). Furthermore since they are immutable they can be used as dictionary keys (lists cannot). To create a tuple the comma-separated list of items has to be enclosed in brackets, or the tuple() operator can be used. Accessing tuple items is done in exactly the same way as lists.

```
atuple = (1, 2, 3)
print ("Length: {}".format(len(atuple)))
print ("First element: {}".format(atuple[0]))
print ("Last element: {}".format(atuple[-1]))

Length: 3
First element: 1
Last element: 3
```

In the next snippet of code it is shown the so called unpacking which is another way to assign tuple values to variables.

```
x, y, z = (10, 5, 12)
print ("coord: x={} y={} z={}".format(x, y, z))

coord: x=10 y=5 z=12
```

If an ntuple has just one element don't forget the comma at the end otherwise it will be treated as a single number.

```
tuple2 = (1,)
print(type(tuple2))
tuple2 = (1)
print(type(tuple2))
```

```
<class 'tuple'>  
<class 'int'>
```

Since a tuple is immutable to add new elements it is necessary to create a new object:

```
tuple1 = (1, 2, 3)  
tuple2 = tuple1 + (4, 5)  
print(tuple2)  
  
(1,2,3,4,5)
```

Finally, as already said tuples can be used as dictionary keys:

```
d = {  
    ('Finance', 1): 'Room 8',  
    ('Finance', 2): 'Room 3',  
    ('Math', 1): 'Room 6',  
    ('Programming', 1): 'IT room'  
}
```

Below the full list of tuple functions:

```
dir(dict)  
  
[...  
    'count',  
    'index']
```

## Chapter 3

# Date and Time

In this chapter we will take a little break and concentrate on a topic that it is not usually covered in this type of courses. However given its importance for financial computation the next paragraphs will be devoted to a close look up on the `datetime` module, whose usage will help in manipulating dates.

### 3.1 Dates

---

As said dates are not usually included in a standard python tutorials, however since they are pretty essential for finance we are going to cover this topic in some detail. In python the date utilities mainly lives in the `datetime` module. We are also going to show `relativedelta` from the `dateutil` module, which allows to add/subtract days/months/years to dates, in other words to make operations on them.

In this first example the today's date is defined and with `relativedelta` two more dates are created adding two months and three days to the first one.

```
from datetime import date, datetime
from dateutil.relativedelta import relativedelta

date1 = date.today()
print (date1)
date2 = date.today() + relativedelta(months=2)
print (date2)
date3 = date.today() - relativedelta(days=3)
print (date3)

2020-08-03
2020-10-03
2020-07-31
```

Here instead another way of computing a new date is shown: in particular a one day delta is stored in a variable and today's date is moved by three days multiplying the defined delta by three.

```
one_day = relativedelta(days=1)
date.today() - 3 * one_day

datetime.date(2020, 7, 31)
```

Next, given two dates their difference is computed (and expressed in days).

```
date1 = date(2019, 7, 2)
date2 = date(2019, 8, 16)
(date2 - date1).days

45
```

Dates can be converted to and from strings and a large variety of formats can be specified in this conversions. The format is determined by a string in which each character starting with % represent an element of the date, e.g. %Y year, %d day, %s seconds, etc...

Below dates to string conversion:

```
date1 = date(2019, 7, 2)
date1.strftime("%Y-%b-%d (%a)") # dates can formatted in many ways
                                # check the docs for more details

'2019-Jul-02 (Tue)'
```

And here, a string is converted to datetime object:

```
# a string can be converted to dates too
datetime.strptime('25 Aug 2019', "%d %b %Y").date()

datetime.date(2019, 8, 25)
```

Finally a last example showing how to get the week-day from a date:

```
date1.weekday() # 0 = Monday, ..., 6 = Sunday

1
```

## Chapter 4

# Python's Object Oriented Programming

In this chapter the main characteristics that makes python an *object oriented programming* language will be reviewed. Before going to OOP however the concepts of function and variable scope will be outlined.

### 4.1 Functions

---

A function is a block of organized, reusable code that is used to perform a single action. Functions provide better modularity for your application and high degree of code reusing. To define a function the keyword `def` is used, followed by the name of the function and by the required parameters in parenthesis. Functions are called by name passing the necessary parameters if any.

```
# sum up all the integers between 1 and n
def my_function(n): # this function take one input only (n)
    x = 0
    for i in range(1, n+1):
        x += i
    return x # the function returns a number

my_function(5) # 5 + 4 + 3 + 2 + 1
```

Functions can return any kind of objects (numbers, strings, lists, complex objects...) but it is not mandatory to have a return value, so you can have functions **without** a return statement (e.g. a function that simply take a string as input and print it to screen with a particular format). In addition the syntax of the return is different from other languages like Visual Basic, the returned object doesn't have to have the same name as the function. Indeed above the variable `x` is returned and not the variable `my_function`. Below an example of function not returning anything.

```
def printing(mystring):
    print((myString).upper())
```

Functions can call other functions (once a function has been defined it can be accessed from everyone within the same file or notebook): here `my_function2` calls `my_function`

```
def my_function_2(n, x):
    return "The result is : {}".format(str(my_function(n)*x))
```

```
my_function_2(5, 10)
```

Functions can also call themselves too (i.e *recursion*). In the next example we will see a function that computes the factorial exploiting the following relationship:

$$\begin{cases} n! = n \times (n-1)! & (\forall n > 1) \\ n! = 1 & (\forall n \leq 1) \end{cases}$$

```
def factorial(n):
    if n <= 1:
        return 1
    else:
        return n * factorial(n-1)

factorial(10)
```

In this example the function `factorial` is initially called with the input corresponding to the factorial we want to compute, it then call itself each time with  $n - 1$ , multiplying together all the results. The previous example is quite simple but recursion can be tricky sometimes so apply it with caution.

Functions input parameters can have default values, which means that a function that works with some input values can be called with less parameters provided their default values have been specified.

In the following example the function `powers` takes three inputs: a list of numbers, an exponent ( $n$ ) and a constant ( $c$ ). The code loops through the provided list of numbers and process them according to the formula  $item^n + c$ , it puts the results in a new list which will be finally returned.

```
def powers(l, n=2, c=0):
    return [item**n+c for item in l]

print (powers([5, 11, 6], 3, 4))
print (powers([5, 11, 6]))

[129, 1335, 220]
[25, 121, 36]
```

As you can see the function is called twice with two different set of parameters: in the first case we pass to it the list of numbers, the exponent and the constant, in the second just the same list of numbers. In the latter case, being defined the default values for  $n$  and  $c$ , the function works perfectly well, fewer inputs are provided and the missing ones are replaced by their defaults.

When calling a function parameters can be passed also by name for clarity, in this case of course the order doesn't matter. Compare the two results below:

```
def func(a, b, c):
    return a + b * c

print (func(c=4, b=2, a=1))
print (func(4, 2, 1))
```



6

In the first case the function is called by name, in the second case the parameter are implicitly assigned according to their position.

Another nice feature of python functions is that we can associate an help message to them so that we can easily check what a function is for by simply asking `help(functionName)`:

```
def powers(l, n=2, c=0):  
    """  
    a shifted power function example  
    """  
    return [item**n+c for item in l]  
  
help(powers)  
  
Help on function powers in module __main__:  
  
powers(l, n=2, c=0)  
    a shifted power function example
```

**Remember, it is always very important to document your code !**

## 4.2 Variable scope

Not all variables are accessible from all parts of our program, and not all variables exist for the entire lifetime of the program. The part of a program where a variable is accessible is called its *scope*.

A variable which is defined in the main body, sometimes referred to as global namespace i.e. the code block which is not indented at all, of a file is called a *global variable*. It will be visible throughout the file, and also inside any file which imports that file. Global variables can have unintended consequences because of their wide-ranging effects, that is why we should almost never use them and they are usually represented by an uppercase name. Only objects which are intended to be used globally, like functions and classes (which will be introduced in the next section), should be put in the global namespace.

Global variables cannot be accessed directly inside functions but before using them they have to be named in a special statement starting with the keyword `global`. Essentially `global` tells python that in the following function we want to use the listed global variable. Imagine a global variable `AGLOBALPARAM` has been defined at the beginning of a program, in the example below few cases are outlined:

- a function that read the value of `AGLOBALPARAM` without modifying it;
- a function that read and modify `AGLOBALPARAM`;
- and a function that throws an exception (i.e. an error in technical language) because it has been badly coded.

```
AGLOBALPARAM = 10
```

```

# Here you just use AGLOBALPARAM value, but do not modify it
# param is just a local copy of AGLOBALPARAM
def multiplyParam(param):
    param = param * 10
    return (param)

# Here you actually use AGLOBALPARAM
# you modify it directly with the global command
def divideParam():
    global AGLOBALPARAM
    AGLOBALPARAM = AGLOBALPARAM / 10
    return (AGLOBALPARAM)

# Here you try to use AGLOBALPARAM but gives you an error
# AGLOBALPARAM is not defined in the function body
# and the global command has not been used neither
def sumParam():
    AGLOBALPARAM = AGLOBALPARAM + 10
    return (AGLOBALPARAM + x)

print ("AGLOBALPARAM is {} to start.".format(AGLOBALPARAM))
print ("Let's multiply it by 10.")
multiplyParam(AGLOBALPARAM)
print ("AGLOBALPARAM is still {}".format(AGLOBALPARAM))
print ("Let's divide it by 10")
divideParam()
print ("Now AGLOBALPARAM is {}".format(AGLOBALPARAM))
print ("Let's sum it to 10")
sumParam()

```

A variable which is defined in a block of code is said to be local to that block. Examples of local scopes are: functions, for or while loops, if blocks, ... In the case of a function it means that a local variable will be accessible from the point at which it is defined until the end of the function itself (e.g. function parameters are examples of local variables).

```

# functions are not evaluated until their are not called
def test_scope(max_val):
    for i in range(max_val):
        print (i)
    print ("max_val in 'test_scope' function is {}".format(max_val))

# the Python interpreter starts evaluating the code from here
max_val = 10
test_scope(5)
print ("max_val in global scope is {}".format(max_val))
print (i)

```

In the previous example we have defined two `max_val` variables, one which is global and it has been initially set to 10, another one which is local to the `test_scope` function. Try as much as

possible to choose different names for each variable you are going to use in a program to avoid confusions and mistakes which may lead to unexpected behaviour of your code.

As a last example compare the following for loops; the first one, correctly written, loops with the variables *i* and *j*, in the second one *j* has been replaced by *i*, note how this is perfectly legal but the result changes dramatically:

```
a = [["a", "b", "c"], ["d", "e", "f"], ["g", "h", "i"]]
for i in range(3):
    for j in range(3):
        print (a[i][j])
```

a  
b  
c  
d  
e  
f  
g  
h  
i

```
a = [["a", "b", "c"], ["d", "e", "f"], ["g", "h", "i"]]
for i in range(3):
    for i in range(3):
        print (a[i][i])
```

a  
e  
i  
a  
e  
i  
a  
e  
i

## 4.3 Classes

Classes are a key ingredient of *Object Oriented Programming* (OOP) and their concept is implemented basically in every modern programming language like python, Java and C++. OOP is a programming model in which programs are organized around data, or objects, rather than functions and logic. **Any object can be thought of a dataset with unique attributes and behaviour** (examples can range from physical entities, such as a human being that is described by properties like name and birthday, down to abstract concepts as a discount curve). This opposes the historical approach to programming where emphasis was placed on how the logic was written rather than how to define the data within the logic. In this framework classes are a mean for creating objects (a particular data structure), providing initial values for its state (member variables or attributes), and implementations of behavior (member functions or methods). Let's summarize here some terminology:

CLASS	→	collection of functions that operate on some dataset
DATA ITEMS	→	are called the <b>attributes</b> of the class
CLASS INSTANCE	→	a specific collection of data belonging to the particular object we are representing
CLASS FUNCTIONS	→	are called <b>methods</b> of the class, and act on <b>instances</b>

**In other words classes are collections of functions that operate on a dataset, and instances of that class represent individual datasets (or in other words a specialization of that class).**

Examples of class are: a class representing a generic building (with number of entrances, number of floors, a flag to know if there is a garden...), a generic dog (with age, fur color...) or a generic computer (with manufacturer, RAM size, CPU type,...). While examples of corresponding instances are: the Empire State Building (a specific building), Lassie (a very particular dog), or your computer, see Fig. 4.1.

To see how they can be defined let's try to code a class representing a person:

```
from datetime import date

class Person:
```

First of all, if needed, we have to import the necessary modules, in this case the `datetime` module is used to managed the person age. Then the `class` keyword followed by the class name is used to start the actual class definition.

### The Constructor Method

After declaring the class name, the constructor method must be defined. In python, this is denoted by `__init__()` regardless the class name. The `__init__` function, as every other method, takes `self` as the first argument, and then any number of arguments as desired by the programmer. The *constructor* allows to specify the initial state of a class by setting its attribute values. For this example that describes a `Person`, the programmer wants to know the name, the birthday and a job (this last one won't be initialize by the constructor).

The `self` parameter is used to create class attributes. Variables whose name starts with `self` have *class scope*, which means are available within each class method. To use the parameters and

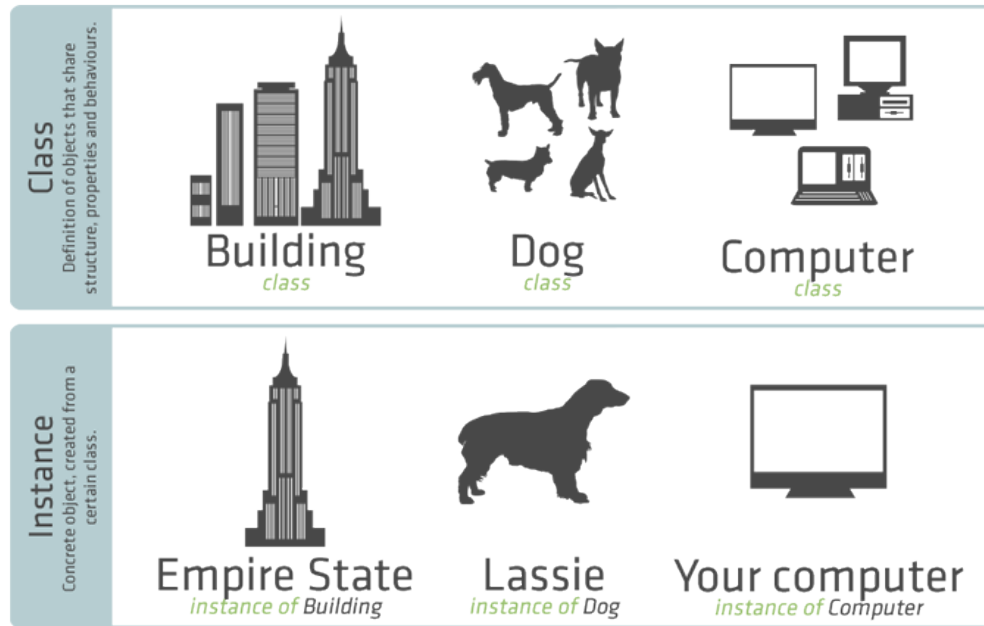


Figure 4.1: Graphical representation of a class instance.

associate them with a particular instance of the class, within the `__init__` method, create variables for each argument like this: `self.variableName = param`.

```
from datetime import date

# this is the class definition
# usually classes use camel naming convention
class Person:
    # the special method __init__ allows to instantiate a class
    # with an initial dataset
    def __init__(self, name, birthday):
        self.name = name
        self.birthday = birthday
        self.occupation = None # this attribute not set at instantiation
```

Now we that have a class definition that represent a generic person we can specialized it to some real person.

```
me = Person("Matteo", date(1974, 10, 20))
print (type(me))

<class '__main__.Person'>
```

Essentially when we instantiate a class python first calls the `__init__` method and initializes the class attributes with the parameter we are passing.

## Class Methods

We haven't yet defined any "person behaviour", so let's add a couple of methods to our class, one computing the person's age and the other setting its primary occupation.

```
class Person:
    def __init__(self, name, birthday):
        self.name = name
        self.birthday = birthday
        self.employment = None

    # this is a normal method and will work on some class attribute
    def age(self, d=date.today()):
        age = (d - self.birthday).days/365
        print("{} is {:.0f} years old".format(self.name, age))

    def mainOccupation(self, occupation):
        self.employment = occupation
        print("{}'s main occupation is: {}".format(self.name, self.employment))
```

To access class attributes and methods the . (dot) operator has to be used.

```
me.name
```

```
'Matteo'
```

```
me.age(date.today())
```

```
Matteo is 46 years old
```

### 4.3.1 Inheritance and Overriding Methods

Inheritance is basically the idea that different classes can have similar components, and in order to avoid repeating code, inheritance is used to link parent classes to descendant classes.

For example, in a fantasy story, there are heroes and monsters but both the heroes and the monsters are characters. And both dragons and orcs are monsters. Though dragons and orcs are different monsters, they share some qualities: they both have a color, they both have a size, they both have enemies. Orcs might have characteristics that dragons do not; for example, what kind of weapon does the orc carry? (Fig. 4.2). Inheritance allows the classes to share information relevant to multiple parts of the code.

Inheritance allows code to be reused and reduces the complexity of a program. The derived classes (descendants) override or extend the functionality of base classes (ancestors). To see how it works in more detail two derived class will be implemented from Person: Adult and Child.

```
class Adult(Person):
    def __init__(self, name, birthday, drv_license_id):
        Person.__init__(self, name, birthday) # this is a special syntax
        self.drv_license_id = drv_license_id
```

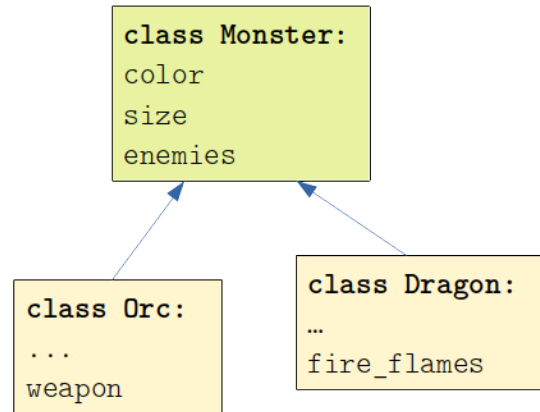


Figure 4.2: Example of inheritance: both Dragon and Orc inherits from Monster, but both add additional attributes that qualify better each "monster features".

```

class Child(Person):
    def mainOccupation(self):
        self.employment = "schoolchild"
        print ("{} is a {}".format(self.name, self.employment))
  
```

Inheritance is specified in parenthesis after the class name. Child still inherits from Person and modifies it by setting the only occupation allowed for a child, overriding the mainOccupation method. Adult class inherits from Person and extends it adding a new attribute (driving license id).

```

pippo = Adult("Goofy", date(1936, 1, 1), "A1234")

qui = Child("Huey", date(2014, 10, 9))
  
```

Behind the scenes, when instantiating a Child class, the constructor of Person is called and the attributes, name and birthday are initialised.

For the Adult class things are a little bit different because a new attribute has been added, so we need to code its own constructor where we first call explicitly the Person's constructor (Person.\_\_init\_\_(self,...) and then we initialise the new attribute.

```

pippo.age()
pippo.mainOccupation("Comic's character")

Goofy is 85 years old
Goofy's main occupation is: Comic's character
  
```

```

qui.age()
qui.mainOccupation()

Huey is 6 years old
Huey is a schoolchild
  
```





## Chapter 5

# Data Manipulation and Its Representation

In this chapter a closer look to a couple more of modules is given. These modules will result to be very useful in managing financial data and to report result of our analysis.

### 5.1 Getting Data

---

The first step of any analysis is usually the one that involves selection and manipulation of data we want to process. Data sources can be various (e.g. website, figures, twitter messages, CSV or Excel files...) and partially reflect its nature which can range from *unstructured* data (without any inherent structure, e.g. social media data) to completely *structured* data (where the data model is defined and usually there is no error associated, e.g. stock trading data).

Our primary goal, before start processing data, is to collect and store the information in a suitable data structure. Python provides a very useful module, called pandas, which allows to collect and save data in *dataframe* objects that can be later on manipulated for analysis purposes.

Looking at pandas manual dataframe are defined as multi-dimensional, size-mutable, potentially heterogeneous, tabular data structure with labeled axes (rows and columns), in much simpler words it is a table whose structure can be modified. It presents data in a way that is suitable for data analysis, contains multiple methods for convenient data filtering and in addition has a lot of utilities to load and save data pretty easily.

Dataframes can be created by:

- importing data from file;
- creating by hand data and then filling the dataframe.

```
import pandas as pd

# reading from file
df1 = pd.read_excel('sample.xlsx') # Excel file
df2 = pd.read_csv('sample.csv') # Comma Separated file

df1.head(11) # show just few rows at the beginning
```

	Date	Price	Volume
0	2000-07-30	100.000000	191.811275
1	2000-07-31	129.216267	190.897541
2	2000-08-01	147.605516	197.476379
3	2000-08-02	107.282251	199.660061
4	2000-08-03	106.036826	200.840459
5	2000-08-04	118.872757	197.130212
6	2000-08-05	101.904544	204.552521
7	2000-08-06	106.392901	198.160030
8	2000-08-06	106.392901	191.125969
9	2000-08-06	106.392901	196.719061
10	2000-08-06	106.392901	196.759837

```
# creating some data in a dictionary
d = {"Nome":["Elisa", "Roberto", "Ciccio", "Topolino", "Gigi"],
      "Età":[1, 27, 25, 24, 31],
      "Punteggio":[100, 120, 95, 1300, 101]}

# filling the dataframe
df = pd.DataFrame(d)
df.head()
```

	Nome	Età	Punteggio
0	Elisa	1	100
1	Roberto	27	120
2	Ciccio	25	95
3	Topolino	24	1300
4	Gigi	31	101

Of course with pandas it is possible to perform a large number of operations on a dataframe. For example it is possible to add a column as a result of an operation on other columns. Looking back at the df1 dataframe it is possible to add a column with the daily variation of the price.

```
import numpy as np

# first let's add an empty column
df1['Variation'] = np.nan # nan stands for not a number

# loop on the Price column, compute the variation and fill the column
# len returns the number of rows of a dataframe
for i in range(1, len(df1)):
    # select the ith row and fill "Variation"
    # loc takes as inputs row and column-name
    df1.loc[i, "Variation"] = (df1.loc[i, "Price"] - df1.loc[i-1, "Price"]) /
                             df1.loc[i-1, "Price"]

df1.head()
```

	Date	Price	Volume	Variation
--	------	-------	--------	-----------

0	2000-07-30	100.000000	191.811275	NaN
1	2000-07-31	129.216267	190.897541	0.292163
2	2000-08-01	147.605516	197.476379	0.142314
3	2000-08-02	107.282251	199.660061	-0.273183
4	2000-08-03	106.036826	200.840459	-0.011609

Of course the first “variation” value is NaN since there is no previous price to compare with.

### 5.1.1 Manage Data

Once we have created our dataframe we may want to preliminary process data to perform very common operations like:

- remove unwanted observations or outliers;
- handle missing data;
- filter, sort and clean data.

### 5.1.2 Unwanted observations and outliers

#### Duplicates

It may happen that our data has duplicates (e.g. those can arise when combining two datasets), or the dataset contains irrelevant fields for the specific study we are carrying on. To find and remove duplicates pandas has convenient methods:

```
# find duplicates based on all columns
# and show just the first 15 results
#print (df1.duplicated()[:15])

# find duplicates based on 'Price'
# and show just the first 15 results
print (df1.duplicated(subset=['Price'])[:15] )
```

```
0    False
1    False
2    False
3    False
4    False
5    False
6    False
7    False
8     True
9     True
10    True
11   False
12   False
13   False
14   False
dtype: bool
```

```
print ("Initial number of rows: {}".format(len(df1)))

# remove duplicates
# where the second argument can be `first`, `last`
# or `False` (consider all of the same values as duplicates).
df1 = df1.drop_duplicates(subset='Price', keep='first')

print ("Number of columns after drop: {}".format(len(df1)))
```

Initial number of rows: 734  
Number of columns after drop: 729

If we would like to drop irrelevant columns for our analysis it is enough to:

```
df2 = df2.drop(columns=['Volume'])
df2.head()
```

	Date	Price
0	2000-07-30	100.000000
1	2000-07-31	129.216267
2	2000-08-01	147.605516
3	2000-08-02	107.282251
4	2000-08-03	106.036826

If instead we just want to remove few rows we can select them by index:

```
# we remove row 0th and 2nd
# axis=0 means use the index column
df2 = df2.drop([0, 2], axis=0)
df2.head()
```

	Date	Price
1	2000-07-31	129.216267
3	2000-08-02	107.282251
4	2000-08-03	106.036826
5	2000-08-04	118.872757
6	2000-08-05	101.904544

Changing the column that act as index we can select the rows also by other attributes:

```
# tell pandas to use Date as index column
df2 = df2.set_index('Date')

# select row to remove by date at this point
df2 = df2.drop(["2000-07-31"], axis=0)

df2.head()
```

Date	Price
2000-08-02	107.282251
2000-08-03	106.036826

```
2000-08-04 118.872757
2000-08-05 101.904544
2000-08-06 106.392901
```

## Outliers

An outlier is an observation that lies outside the overall pattern of a distribution. Common causes can be human, measurement or experimental errors. Outliers must be handled carefully and we should remove them cautiously, *outliers are innocent until proven guilty*. We may have removed the most interesting part of our dataset !

The core statistics about a particular column can be studied by the `describe()` method which returns the following information:

- for numeric columns: the value count, mean, standard deviation, minimum, maximum and 25th, 50th and 75h quantiles for the data in a column;
- for string columns: the number of unique entries, the most frequent occurring value (*top*), and the number of times the top value occurs (*freq*).

```
df1.describe()

      Price      Volume  Variation
count  728.000000  729.000000  724.000000
mean   120.898678  200.355900    0.146330
std    490.493411    4.970745    3.637952
min      0.878873  186.430551   -0.995284
25%    14.809934  196.998603   -0.119423
50%    61.325699  200.221125   -0.005549
75%   164.021813  203.580691    0.121290
max  13000.000000  215.140868   97.756432
```

Looking at mean and std and comparing it with min and max values we could find a range outside of which we may have outliers. For example 13000.0 is several standard deviation away the mean which may indicate that it is not a good value.

Another way to spot outliers is to plot column distributions and again pandas comes to help us:

```
df1.hist("Variation", bins=np.arange(0, 100, 1))

-----

NameError                                Traceback (most recent call last)

<ipython-input-1-97dbdc6fcfec> in <module>
----> 1 df1.hist("Variation", bins=np.arange(0, 100, 1))

NameError: name 'df1' is not defined
```

From the histograms it is clear how the value of 97.76, is far from general population. This doesn't mean they are necessarily wrong but it should make ring a bell in our head...

To remove outliers from data we can either remove the entire rows or replace the suspicious values by a default value (e.g. 0, 1, a threshold value...).

**Note:** missing data may be informative itself ! When filling the gap with *artificial data* (e.g. mean, median, std...) having similar properties than real observation, the added value won't be scientifically valid, no matter how sophisticated your filling method is.

```
import numpy as np
```

```
df2.replace(1300, 500)      # replace 1300 with 500
df2 = df2.replace(1300, np.nan)  # replace 1300 with NaN

df2 = df2.mask(df1 >= 600, 500)  # replace every element >=600 with 5
```

### 5.1.3 Handle Missing Data

Usually when importing data with pandas we may have some NaN values (short for *not a number* which represent the null value). NaN is the value that is given to missing fields in a row. Like for the outliers we can use the replace or mask methods to remove the NaNs. In case the whole row as NaN it may be wise to drop it entirely.

Additionally we can use dropna() which remove all the NaN at once.

```
df1 = df1.dropna()

print ("Number of rows after dropping NaN: {}".format(len(df1)))

Number of rows after dropping NaN: 724
```

### 5.1.4 Filter, Sort and Clean Data

#### Filtering

When we work with huge datasets we may reach computational limits (e.g. insufficient memory, CPU performance, too slow processing time...) and in those cases it can be helpful to filter data by attributes for example by splitting by time or some other property.

Assuming to have the following table and putting back the volume column

```
# df.iloc[row, col]
# NOTE: iloc takes row and column index (two numbers)
# loc instead takes row index and column name
print (df1.iloc[1, 2]) # returns 62 the volume associated with the row 1

print()
#df.iloc[row1:row2, col1:col2]
# this is called slicing, remember ?
print (df1.iloc[0:2, 2:3]) # returns rows 0 and 1 of column 2

197.476378531652

      Volume
1  190.897541
```

```
2 197.476379
```

```
subset = df1.iloc[:, 1] # select column 1

subset = df1.iloc[2, :] # select row 2

subset = df1.iloc[0:2, :] # select 2 rows

subset = df1.iloc[:, 2, :] # this is equivalent to before
```

A more advanced way of filtering is the following (it apply a selection on the values). The notation is a bit awkward but very useful:

```
import datetime

# colon means all the rows
subset = df1[df1.iloc[:, 0] < datetime.datetime(2000, 8, 15)]
print(subset)
```

	Date	Price	Volume	Variation
1	2000-07-31	129.216267	190.897541	0.292163
2	2000-08-01	147.605516	197.476379	0.142314
3	2000-08-02	107.282251	199.660061	-0.273183
4	2000-08-03	106.036826	200.840459	-0.011609
5	2000-08-04	118.872757	197.130212	0.121052
6	2000-08-05	101.904544	204.552521	-0.142743
7	2000-08-06	106.392901	198.160030	0.044045
11	2000-08-07	107.646053	198.861429	0.011779
12	2000-08-08	106.666468	197.213497	-0.009100
13	2000-08-09	101.981029	204.425797	-0.043926
14	2000-08-10	110.100330	196.122844	0.079616
15	2000-08-11	138.656481	200.703360	0.259365
16	2000-08-12	113.180782	205.676449	-0.183732
17	2000-08-13	137.639947	203.468517	0.216107
18	2000-08-14	142.646169	198.528626	0.036372

## Sorting

To sort our data we can use `sort_values()` method (it can be specified ascending, descending).

```
# sort by price then by date in descending order
df2.sort_values(by=['Price', "Date"], ascending=False)[:10]
```

Date	Price
2000-08-20	13000.000000
2000-10-20	593.477666
2001-01-05	571.444679
2000-12-31	532.558487
2000-10-14	516.044122
2001-01-02	503.583189

2001-01-01	502.849987
2000-12-30	487.353466
2001-01-04	478.027182
2001-01-10	473.061993

### Cleaning or Regularizing

As we will see when dealing with machine learning, often we need to regularize our data to improve the stability of a training. One typical situation is when we want to *normalize* data, which means re-scale the values into a range of [0, 1].

$$x = [1, 43, 65, 23, 4, 57, 87, 45, 45, 23]$$

$$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

$$x_{new} = [0, 0.48, 0.74, 0.25, 0.03, 0.65, 1, 0.51, 0.51, 0.25]$$

To apply such a transformation with pandas is very easy since applying the formula to a dataframe implies it is done to each row:

```
df1['Price'] = (df1['Price'] - df1['Price'].min()) \
    / (df1['Price'].max() - df1['Price'].min())
df1.head()
```

	Date	Price	Volume	Variation
1	2000-07-31	0.009873	190.897541	0.292163
2	2000-08-01	0.011287	197.476379	0.142314
3	2000-08-02	0.008185	199.660061	-0.273183
4	2000-08-03	0.008090	200.840459	-0.011609
5	2000-08-04	0.009077	197.130212	0.121052

Another quite common transformation is called *standardization*, essentially we re-scale data to have 0 mean and standard deviation of 1:

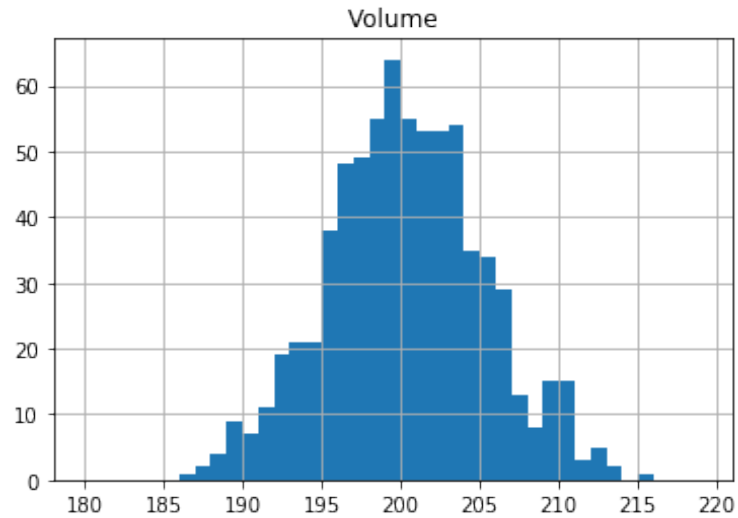
$$x_{new} = \frac{x - \mu}{\sigma}$$

Again it is straightforward to do it in pandas:

```
df1.hist('Volume', bins=np.arange(180, 220, 1))
print (df1['Volume'].mean())
print (df1['Volume'].std())

200.36750575214748
4.968224698257929
```

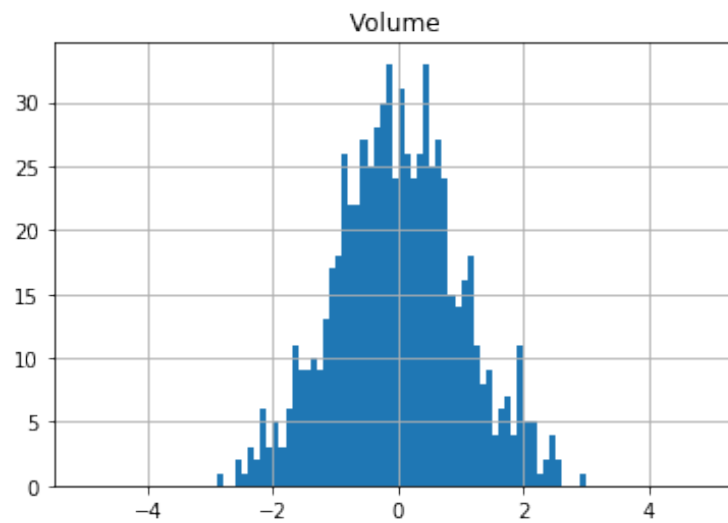




```
df1['Volume'] = (df1['Volume'] - df1['Volume'].mean()) / df1['Volume'].std()

df1.hist('Volume', bins=np.arange(-5, 5, 0.1))
print (df1['Volume'].mean())
print (df1['Volume'].std())

-6.148550054609154e-15
1.0
```



## 5.2 Plotting in python

As we have just seen pandas allows to quickly draw histograms of dataframe columns, but during an analysis we may want to plot distributions from list or objects not stored in a dataframe. Furthermore the simple and very useful provided interface doesn't grant full access to all histogram features that we need to produce nice and informative plots.

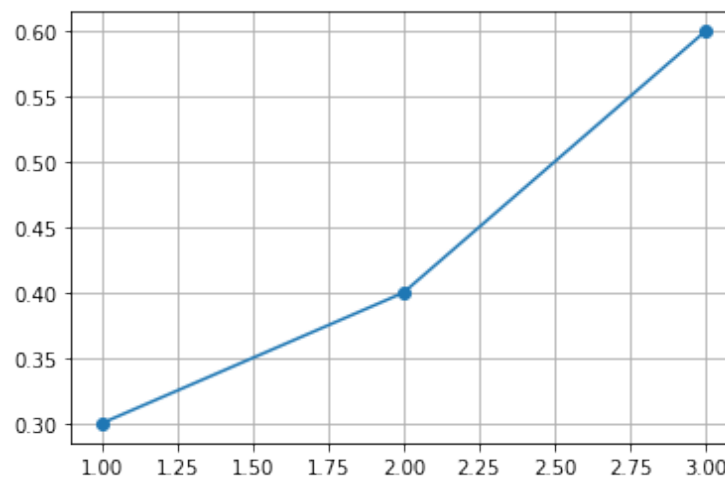
In order to do so we can use the `matplotlib` module which is specifically dedicated to plotting (pandas interface is based on the same module indeed). Let's look briefly to its capability by examples.

### 5.2.1 Plot a graph given $x$ and $y$ values (scatter-plot)

```
from matplotlib import pyplot as plt

x = [1, 2, 3]
y = [0.3, 0.4, 0.6]

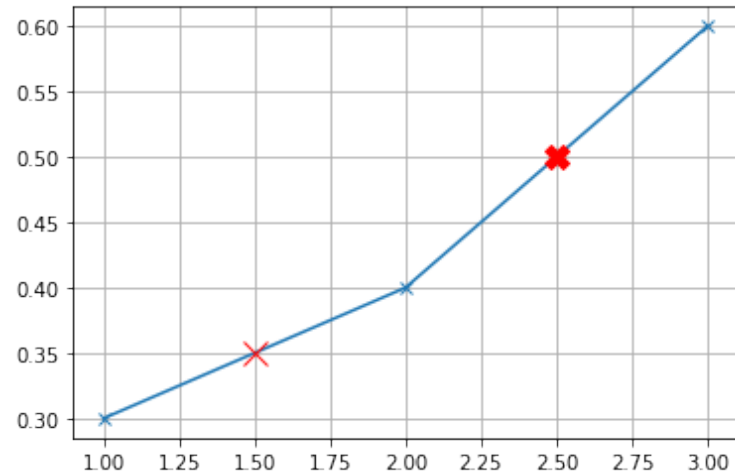
plt.plot(x, y, marker='o') # we are using circle markers
plt.grid(True)             # this line activate grid drawing
plt.show()
```



```
# if we want to plot specific points too

x = [1, 2, 3]
y = [0.3, 0.4, 0.6]

plt.plot(x, y, marker='x')
plt.plot(2.5, 0.5, marker='X', ms=12, color='red')
plt.plot(1.5, 0.35, marker='x', ms=12, color='red')
plt.grid(True)
plt.show()
```

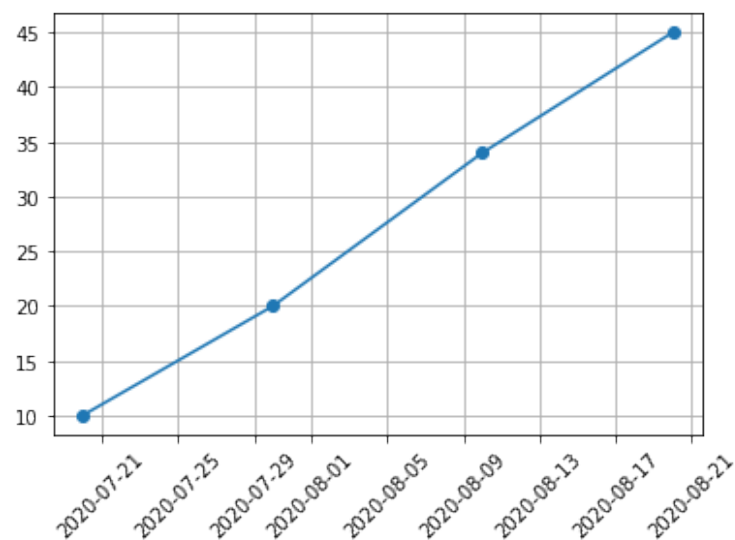


What if  $x$  values are dates ?

```
import datetime
import matplotlib.dates as mdates

x = [datetime.date(2020, 7, 20), datetime.date(2020, 7, 30),
      datetime.date(2020, 8, 10), datetime.date(2020, 8, 20)]

y = [10, 20, 34, 45]
plt.plot(x, y, marker='o')
# this line tells matplotlib we have dates on x axis
plt.gca().xaxis.set_major_formatter(mdates.DateFormatter('%Y-%m-%d'))
# this one instead rotate labels to avoid superimposition
plt.xticks(rotation=45)
plt.grid(True)
plt.show()
```

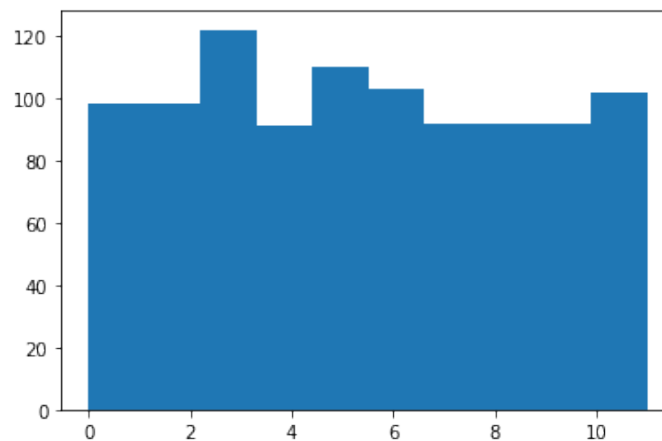


### 5.2.2 Plotting an Histogram

```
import random
numbers = []
for _ in range(1000):
    numbers.append(random.randint(1, 10))

from matplotlib import pyplot as plt

# Here we define the binning
# 6 is the number of bins, going from 0 to 10
plt.hist(numbers, 10, range=[0, 11])
plt.show()
```



### Plotting a Function

In this case let's try to make the plot prettier adding labels, legend... All the commands apply also to the previous examples.

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import norm

# define the functions to plot
# a gaussian with mean=0 and sigma=1
# in scipy module this is called norm
mu=0
sigma = 1
x = np.arange(-10, -1.645, 0.001)
x_all = np.arange(-4, 4, 0.001)
y = norm.pdf(x, 0, 1)
y_all = norm.pdf(x_all, 0, 1)

# draw the gaussian
```

```

plt.plot(x_all, y_all, label='Gaussian')

# fill with different alpha using x_all and y_all as limits
# alpha set the transparency level: 0 transparent, 1 solid
plt.fill_between(x_all, y_all, 0, alpha=0.1, color='blue', label="Gaussian CDF")

# fill with color red using x and y as limits
# label associate text to the object for the legend
plt.fill_between(x, y, 0, alpha=1, color='red', label="5% tail")

# set x axis limits
plt.xlim([-4, 4])

# add a label for X axis
plt.xlabel("Changes of value")

# add a label to y axis
plt.ylabel("Gaussian values")

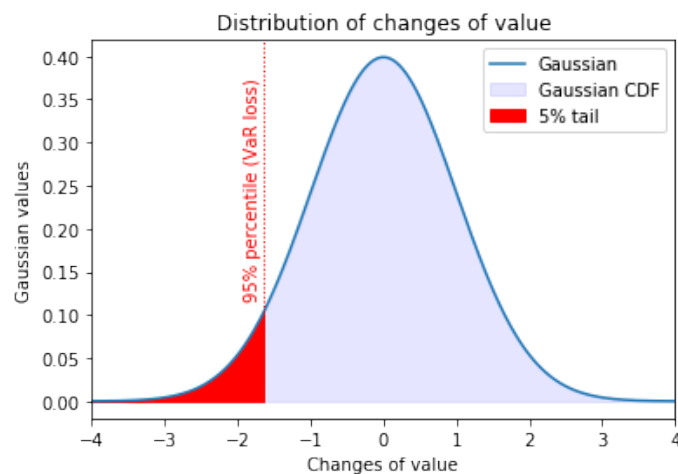
# add histogram title
plt.title("Distribution of changes of value")

# draw a vertical line at x=-1.645
# y limits are in percent w.r.t. to y axis length
plt.axvline(x=-1.645, ymin=0.1, ymax=1, linestyle=':', linewidth=1, color = 'red')

# write some text to explain the line
plt.text(-1.9, .12, '95% percentile (VaR loss)', fontsize=10, rotation=90,
        color='red')

plt.legend()
plt.show()

```



If you are particularly satisfied by your work you can save the graph to a file:

```
plt.savefig('normal_curve.png')
```

```
<Figure size 432x288 with 0 Axes>
```

## Chapter 6

# Interpolation, Discount Factors and Forward Rates

In this chapter we will start to see the first applications of python to financial calculations. In particular we will consider discount curves and forward rates, implementing the first utilities that will fill our financial module. In doing so we will review a widely used mathematical tool: *interpolation*.

### 6.1 Linear interpolation

---

Consider to have few data points, obtained by sampling or experimenting. These points represent the values of a not well known function  $f(x)$ , where  $x$  is an independent variable (e.g. in recording a trip: distances at certain times,  $d = f(t)$ ).

It may be necessary to estimate values of the function  $f$  at values for which we don't have samples. Interpolation is a method of "constructing" new points within the range of the known data.

Let's clarify the technique with an example. Assume you are going on holidays by car and that luckily there isn't much traffic so that you can drive at constant speed (which gives a linear relation between traveled space and time i.e.  $s = v \cdot t$ , which means that if you plot the distances  $s$  as a function of the time  $t$  you get a line with slope  $v$ ), see Fig. 6.1.

Given two samples of the car traveled distance  $s_1$  and  $s_2$  taken at two different times  $t_1$  and  $t_2$  you can linearly interpolate to find your position at different times using the following relations:

$$s = (1 - w) \cdot s_1 + w \cdot s_2$$

where  $t$  is a generic time at which we want to know the distance  $s$  and  $w = \frac{t - t_1}{t_2 - t_1}$ .

#### Derivation

The equation of a line for two points  $(t_1, s_1)$  and  $(t_2, s_2)$  can be written as:

$$\frac{t - t_1}{t_2 - t_1} = \frac{s - s_1}{s_2 - s_1}$$

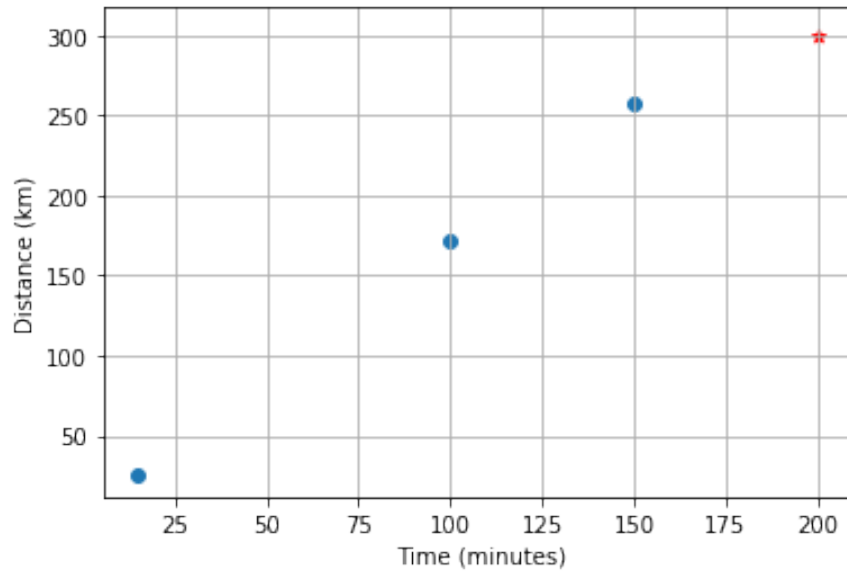


Figure 6.1: An example of sampling of traveled distances at some time. The red point shows an additional sample taken after the trip velocity has been reduced.

Setting again  $w = \frac{t - t_1}{t_2 - t_1}$  and solving for  $s$  we find the desired solution:

$$(s_2 - s_1) \cdot w = s - s_1 \quad \Rightarrow \quad s = (1 - w) \cdot s_1 + w \cdot s_2$$

This formula can also be understood as a weighted average where the weights are inversely related to the distance from the end points to the unknown point ( $w_1 = (1 - w) = \frac{t_2 - t}{t_2 - t_1}$ ,  $w_2 = w$ ), the closer point has more influence than the farther point.

Back to our example, if  $s_1 = 25.75$  km (@ $t_1 = 15$  min) and  $s_2 = 171.7$  km (@ $t_2 = 100$  min) let's find distance traveled in 1 hour (interpolation):

```
s_1 = 25.75 # distance in km
t_1 = 15    # elapsed time in minutes
s_2 = 171.7
t_2 = 100

t = 60

w = (t - t_1)/(t_2 - t_1)
s = (1 - w)*s_1 + w*s_2

print("{:.1f} km".format(s))

103.0 km
```

Always interpret critically your results to guess if they make sense or not. In the previous example we certainly expected something between 25.75 and 171.7 km (our range ends) furthermore since we are looking for the distance at a time which is almost halfway the interval, the result will



be somehow in the middle or around 98.6 km. This is indeed more or less what we have got. This simple reasoning should be applied every time you have a result to quickly judge it.

If we believe the relation between our variable stays the same ( $f(t)$  still linear), we can use the same formula to *extrapolate* values *outside* our initial sample. For example if we keep the same constant velocity in our trip we could check the distance traveled after 3 hours:

```
s_1 = 25.75 # distance in km
t_1 = 15    # elapsed time in minutes
s_2 = 171.7
t_2 = 100

t = 180

w = (t - t_1)/(t_2 - t_1)
s = (1 - w)*s_1 + w*s_2

print("{:.1f} km".format(s))

309.1 km
```

### 6.1.1 Log-linear interpolation

When the function  $f$  that we want to interpolate is an exponential we can fall back to the previous case by a simple variable transformation. Assume the following is the relationship between  $p$  and  $h$ , two generic variables:

$$p = \exp(c \cdot h)$$

Applying the logarithm to both sides of the equation gives:

$$s = \log(p) = \log(\exp(c \cdot h)) = c \cdot h$$

so there is linear relation between the new variable  $s$  and  $h$ . At this point we can use the results of the previous section to interpolate for values of  $s$ , just remember to exponentiate the result to get the correct  $p$ . In formulas:

$$w = \frac{h - h_1}{h_2 - h_1}$$

$$s = (1 - w) \cdot s_1 + w \cdot s_2 \quad (\text{remember now } s = \log(p))$$

$$p = \exp(s)$$

Let's see another example. Atmospheric pressure decreases with the altitude (i.e. the highest you flight the lower is the pressure) following an exponential law:

$$p = p_0 \cdot e^{-\alpha h}$$

where

- $h$  is the altitude

- $p_0$  is the pressure at sea level
- $\alpha$  is a constant

Taking the logarithm of each side of the equation I get a linear relation which can be interpolated as seen before:

$$s = \log(p) = \log(p_0 \cdot e^{-\alpha h}) \propto -\alpha \cdot h$$

Now assume that we have measured  $p_1 = 90$  kPa ( $h_1 = 1000$  m) and  $p_2 = 40$  kPa ( $h_2 = 7000$  m) what will be the atmospheric pressure on top of the Mont Blanc (4812 m) ? and on top of Mount Everest (8848 m) ?

```
# pressure on top of the Mont Blanc (interpolation)
from math import log, exp

# first we take the logarithm of our measurements to use the linear
# relation to interpolate
h_1 = 1000 # height in meters
s_1 = log(90) # logarithm of the pressure at height h1
h_2 = 7000 # height in meters
s_2 = log(40) # logarithm of the pressure at height h2

h = 4812

w = (h - h_1)/(h_2 - h_1)
s = (1 - w)*s_1 + w*s_2

print("{:.1f} kPa".format(exp(s)))

53.8 kPa
```

```
# pressure on top of the Mount Everest (extrapolation)
from math import log, exp

# first we take the logarithm of our measurements to use the linear
# relation to interpolate
h_1 = 1000 # height in meters
s_1 = log(90) # logarithm of the pressure at height h1
h_2 = 7000 # height in meters
s_2 = log(40) # logarithm of the pressure at height h2
h = 8848

w = (h - h_1)/(h_2 - h_1)
s = (1 - w)*s_1 + w*s_2

print("{:.1f} kPa".format(exp(s)))

31.2 kPa
```

In this case we check our results by plotting the found pressures on top of the  $P$  vs  $h$  plot shown on Wikipedia, see Fig 6.2.

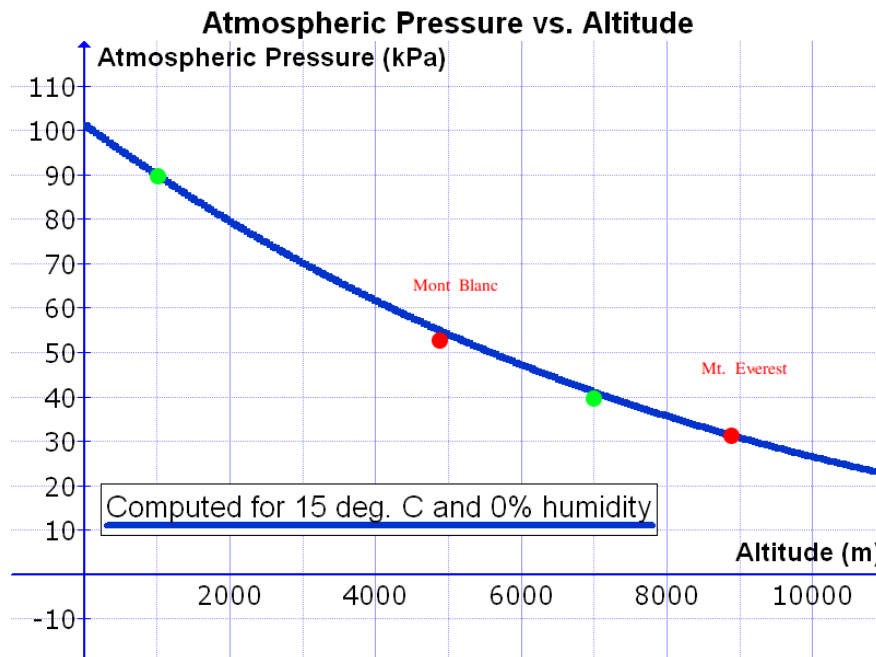


Figure 6.2: Atmospheric pressure versus altitude (Wikipedia). Green points represent our measurements, red points represent interpolation/extrapolation.

### 6.1.2 Limitations of Interpolation

Interpolation is just an approximation and works well when either the function  $f$  is linear or we are trying to interpolate between two points that are close enough to believe that  $f$  is almost linear in that interval.

It can be easily demonstrated that the linear approximation between two points of a given function  $f(x)$  gets worse with the second derivative of the function that is approximated ( $f''(x)$ ). This is intuitively correct: the "curvier" the function is, the worse the approximation made with simple linear interpolation becomes, see Fig. 6.3 where we try to interpolate a sine function.

To improve the approximation accuracy with complicated curves a polynomial of higher order can be used ( $p(x) = a_0 + a_1x + a_2x^2 + \dots$ ), for example in the evaluation of the natural logarithm and trigonometric functions. It has to be clear however that going to higher degrees does not always help (for those interested see [Runge's phenomenon](#)).

## 6.2 Discount curve interpolation

Finally we can come back to finance. Since discount factors are derived from a discrete set of dates we may need to find the factor at some different date and clearly we can use interpolation to do it. Now we will see how to implement a python function which interpolates some given discount factors. Needed data:

- a list of pillars dates specifying the value dates of the given discount factors,  $t_0, \dots, t_{n-1}$ ;
- a list of given discount factors,  $D(t_0), \dots, D(t_{n-1})$ ;
- a pricing date ('today' date) which corresponds to  $t = 0$ .

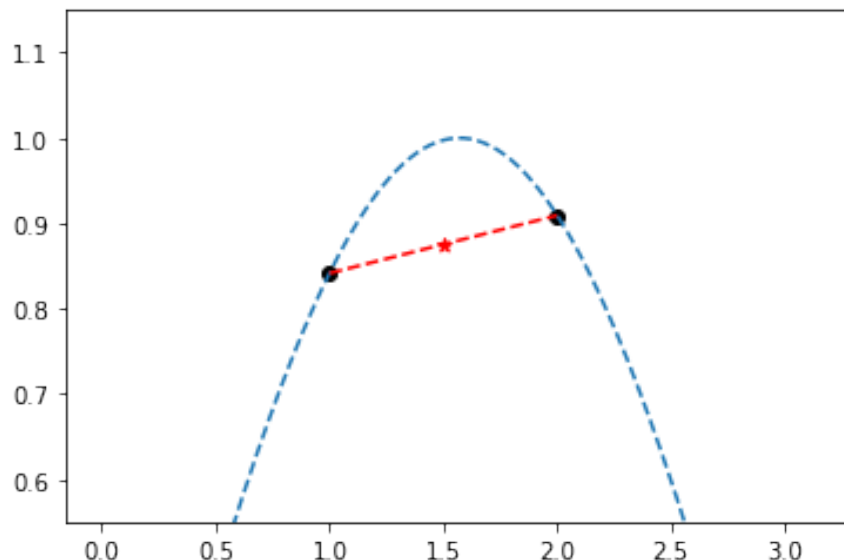


Figure 6.3: Trying to approximate a sine function with a line is clearly not going to work unless the interpolation interval is very small.

The input argument to the function will be the value date at which we want to interpolate the discount factor. Since the discount factor can be expressed as  $D = e^{-r(T-t)}$  the function will use a log-linear interpolation to return the value at a date not included in the given pillars. More technically we can say that we are doing a linear interpolation over time in the log space:

$$d(t_i) := \ln(D(t_i))$$

$$d(t) = (1 - w)d(t_i) + wd(t_{i+1}); \quad w = \frac{t - t_i}{t_{i+1} - t_i}$$

$$D(t) = \exp(d(t))$$

where again  $i$  is such that  $t_i \leq t \leq t_{i+1}$

This time instead of reinventing the wheel and performing the interpolation with our own code, we'll use the function `interp` provided by the module `numpy`; this function linearly interpolates between the provided points to estimate the value of  $f$  at some "new"  $x$ . Say we want to interpolate the points at  $x = 2.5$  given the following values:

```
import numpy as np

xp = [0, 1, 5]
fp = [0, 2, 4]
np.interp(2.5, xp, fp)

2.75
```

Assume we have three discount factors instead:

```
# import modules and objects that we need
from datetime import date
```

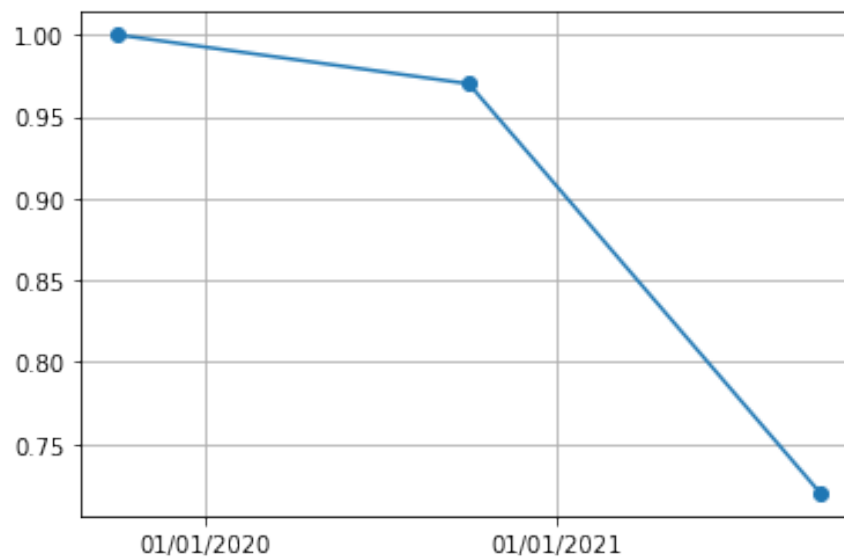
```
import numpy, math
from matplotlib import pyplot as plt
import matplotlib.dates as mdates
# with this notation we tell python to use mdates as an alias
# for matplotlib.dates

# define the input data
today_date = date(2019, 10, 1)

pillar_dates = [date(2019, 10, 1), date(2020, 10, 1), date(2021, 10, 1)]
discount_factors = [1.0, 0.97, 0.72]
```

Let's see what this fake discount curve looks like when plotted on a graph:

```
plt.plot(pillar_dates, discount_factors, marker='o')
plt.gca().xaxis.set_major_formatter(mdates.DateFormatter('%m/%d/%Y'))
plt.gca().xaxis.set_major_locator(mdates.YearLocator())
plt.grid(True)
plt.show()
```



Since it is a computation that from now on we need to perform quite often it is convenient to write a function that compute the discount factor at arbitrary dates.

```
# define the df function
def df(d, pillar_dates, discount_factors):
    # first thing we need to do is to apply the logarithm function
    # to the discount factors since we are doing log-linear and
    # not just linear interpolation
    log_discount_factors = []
    for discount_factor in discount_factors:
        log_discount_factors.append(math.log(discount_factor))
```

```
# perform the linear interpolation of the log discount factors
interpolated_log_discount_factor = \
    numpy.interp(d, pillar_dates, log_discount_factors)

# return the interpolated discount factor
return math.exp(interpolated_log_discount_factor)
```

This is almost OK, **but it won't work** because `numpy.interp` only accepts numbers or a list of numbers as argument i.e. it doesn't automatically convert or interpret dates as numbers so doesn't know how to interpolate them. So we need to do the conversion ourselves before passing the dates into the interpolation function. The following updated version of our function converts the pillar dates into "pillar days" i.e. each date is replaced by the number of days today ( $t_0$ ):

```
def df(d, today_date, pillar_dates, discount_factors):
    # first thing we need to do is to apply the logarithm function
    # to the discount factors since we are doing log-linear and
    # not just linear interpolation
    log_discount_factors = []
    for discount_factor in discount_factors:
        log_discount_factors.append(math.log(discount_factor))

    # convert the pillar dates to pillar 'days'
    # i.e. number of days from today
    # to write shorter code we can use this NEW notation
    # which condenses for and list creation in one line
    pillar_days = \
        [(pillar_date - today_date).days for pillar_date in pillar_dates]

    # obviously we need to do the same to the value date
    # argument of the df function
    d_days = (d - today_date).days

    # perform the linear interpolation of the log discount factors
    interpolated_log_discount_factor = \
        numpy.interp(d_days, pillar_days, log_discount_factors)

    # return the interpolated discount factor
    return math.exp(interpolated_log_discount_factor)
```

Now we can use the `df` function to get discount factors on value dates between the given pillar dates:

```
d0 = date(2020, 1, 1)
df0 = df(d0, today_date, pillar_dates, discount_factors)
print (df0)

0.9923728228571693
```

```
d1 = date(2021, 1, 1)
df1 = df(d1, today_date, pillar_dates, discount_factors)
```

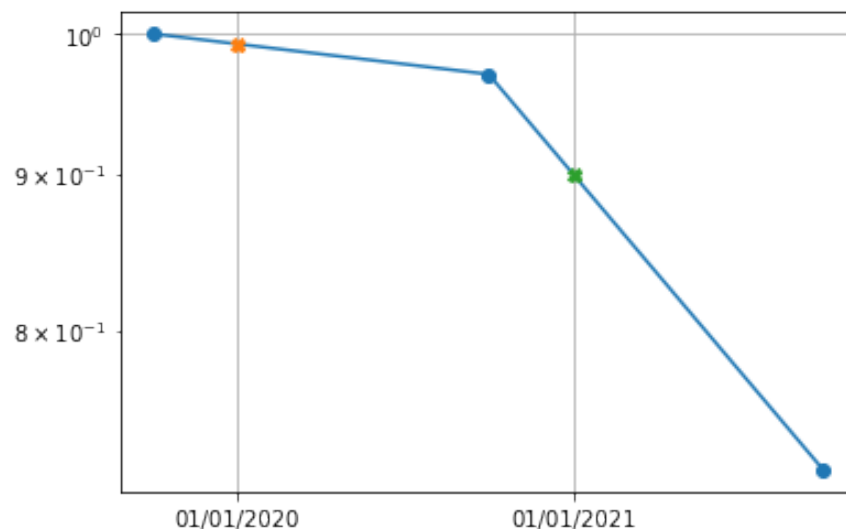
```
print (df1)

0.8997999273630835
```

Another very useful way to check the correctness of a result is by plotting data, so let's see what these look like when plotted on a semi-log graph and if they make sense:

```
from matplotlib import pyplot as plt
import matplotlib.dates as mdates

plt.semilogy(pillar_dates, discount_factors, marker='o')
plt.semilogy(d0,df0 , marker='X')
plt.semilogy(d1,df1 , marker='X')
plt.gca().xaxis.set_major_formatter(mdates.DateFormatter('%m/%d/%Y'))
plt.gca().xaxis.set_major_locator(mdates.YearLocator())
plt.grid(True)
plt.show()
```

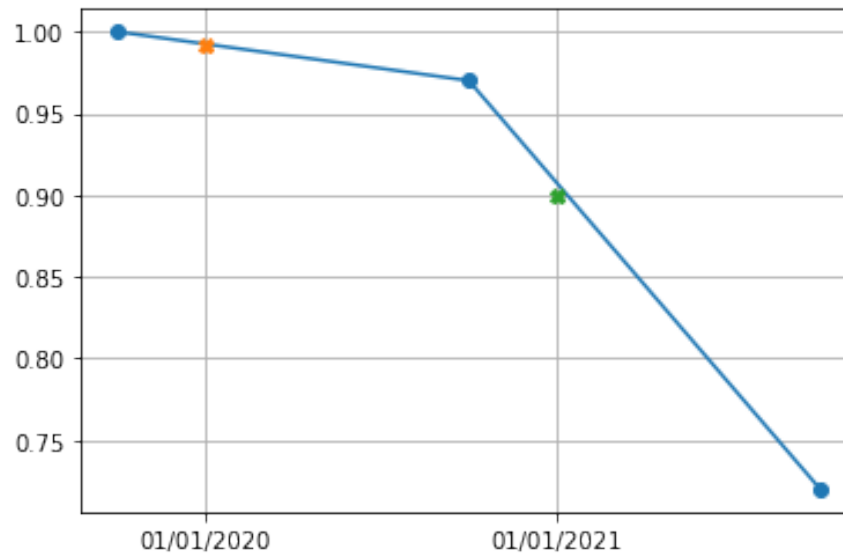


Let's see what these look like when plotted on a linear graph instead:

```
from matplotlib import pyplot as plt
import matplotlib.dates as mdates

plt.plot(pillar_dates, discount_factors, marker='o')
plt.plot(d0,df0 , marker='X')
plt.plot(d1,df1 , marker='X')
plt.gca().xaxis.set_major_formatter(mdates.DateFormatter('%m/%d/%Y'))
plt.gca().xaxis.set_major_locator(mdates.YearLocator())
plt.grid(True)
plt.show()
```

Discrepancies in the linear plot are most likely due to rounding.



### 6.3 Forward Rates

A forward rate is an interest rate applicable to a financial transaction that will take place in the future. Forward rates are calculated from the spot rate by exploiting the no arbitrage condition which states that investing at rate  $r_1$  for the period  $(0, T_1)$  and then *re-investing* at rate  $r_{1,2}$  for the time period  $(T_1, T_2)$  is equivalent to invest at rate  $r_2$  for the full time period  $(0, T_2)$ . Essentially two investors shouldn't be able to earn money from arbitraging between different interest periods. That said:

$$(1 + r_1 T_1)(1 + r_{1,2}(T_2 - T_1)) = 1 + r_2 T_2$$

Solving for  $r_{1,2}$  leads to

$$F(T_1, T_2) = r_{1,2} = \frac{1}{T_2 - T_1} \left( \frac{D(T_1)}{D(T_2)} - 1 \right) \quad (\text{where } D(T_i) = \frac{1}{1 + r_i T_i})$$

```
from datetime import date
import numpy, math

today_date = date(2019, 1, 1)

pillar_dates = [date(2019, 1, 1),
                 date(2020, 1, 1),
                 date(2021, 10, 1)]
discount_factors = [1.0, 0.97, 0.72]

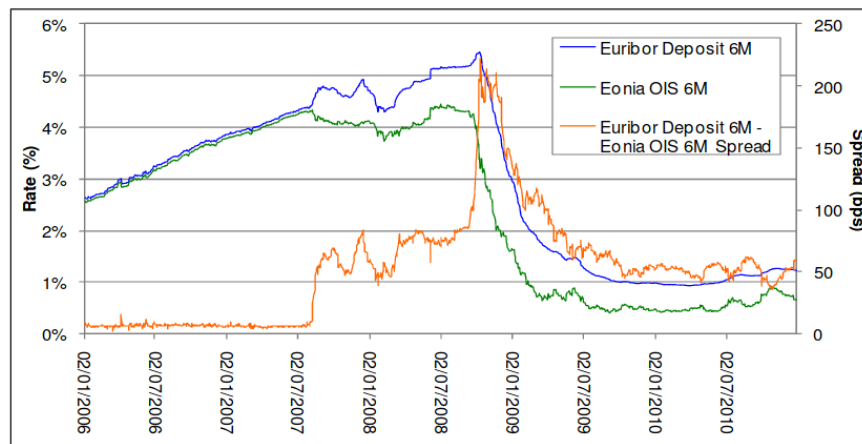
def forward_rate(t1, t2, today_date, pillar_dates, discount_factors):
    return 365.0/(t2-t1).days *
        (df(t1, today_date, pillar_dates, discount_factors) /
         df(t2, today_date, pillar_dates, discount_factors) - 1)
```



```
forward_rate(date(2019, 2, 1), date(2019, 8, 1),
             today_date, pillar_dates, discount_factors)
```

### 6.3.1 2008 Financial Crisis

Looking at the historical series of the Euribor (6M) rate versus the Eonia Overnight Indexed Swap (OIS-6M) rate over the time interval 2006-2011 it becomes apparent how before August 2007 the two rates display strictly overlapping trends differing of no more than 6 bps.



**Figure 1:** historical series of Euribor Deposit 6M rate versus Eonia OIS 6M rate. The corresponding spread is shown on the right axis (Jan. 06 – Dec. 10 window, source: Bloomberg).

In August 2007 however we observe a sudden increase of the Euribor rate and a simultaneous decrease of the OIS rate that leads to the explosion of the corresponding basis spread, touching the peak of 222 bps in October 2008, when Lehman Brothers filed for bankruptcy. Successively the basis has sensibly reduced and stabilized between 40 bps and 60 bps (notice that the pre-crisis level has never been recovered). The same effect is observed for other similar couples of series, e.g. Euribor 3M vs OIS 3M.

The reason of the abrupt divergence between the Euribor and OIS rates can be explained by considering both the monetary policy decisions adopted by international authorities in response to the financial turmoil, and the impact of the credit crunch on both credit and liquidity risk perception of the market, coupled with the different financial meaning and dynamics of these rates.

- The Euribor rate is the reference rate for over-the-counter (OTC) transactions in the Euro area. It is defined as the rate at which Euro inter-bank deposits are being offered within the EMU zone by one prime bank to another at 11:00 a.m. Brussels time. The rate fixings for a strip of 15 maturities (from one day to one year) are constructed as the average of the rates submitted (excluding the highest and lowest 15% tails) by a panel of 42 banks, selected among the EU banks with the highest volume of business in the Euro zone money markets, plus some large international bank from non-EU countries with important euro zone operations. *Thus, Euribor rates reflect the average cost of funding of banks in the inter bank market at each given maturity. During the crisis the solvency and solidity of the whole financial sector was brought into question and the credit and liquidity risk and uremia associated to inter-bank counter-parties sharply increased.* The Euribor rates immediately reflected these dynamics and raise to their highest values over more than 10 years. As seen in the plot above, the Euribor 6M rate suddenly increased on August 2007 and reached 5.49% on 10th October 2008.

- The Eonia rate is the reference rate for overnight OTC transactions in the Euro area. It is constructed as the average rate of the overnight transactions (one day maturity deposits) executed during a given business day by a panel of banks on the inter-bank money market, weighted with the corresponding transaction volumes. *The Eonia Contribution Panel coincides with the Euribor Contribution Panel, thus Eonia rate includes information on the short term (overnight) liquidity expectations of banks in the Euro money market. It is also used by the European Central Bank (ECB) as a method of effecting and observing the transmission of its monetary policy actions. During the crisis the central banks were mainly concerned about stabilizing the level of liquidity in the market, thus they reduced the level of the official rates.* Furthermore, the daily tenor of the Eonia rate makes negligible the credit and liquidity risks reflected on it: for this reason the OIS rates are considered the best proxies available in the market for the risk-free rate.

As a practical result of the divergence of the two indices, after the 2008 financial crisis, it is not possible anymore to use a single discount curve to correctly price forward rates of all tenors. For example, if we want to calculate the net present value of a forward 6-month Libor coupon, we need to simultaneously use two different discount curves:

- the 6-month Libor curve for determining the forward rate;
- the EONIA curve for discounting the expected cash flow.

Our financial library will have to implement the following calculation:

$$\text{NPV} = D_{\text{EONIA}}(T_1) \times \frac{1}{T_2 - T_1} \left( \frac{D_{\text{LIBOR}}(T_1)}{D_{\text{LIBOR}}(T_2)} - 1 \right)$$

and this will be asked to be done in the exercises relative to this chapter.

To exploit the Object Oriented capabilities we will implement a `DiscountCurve` class, below a skeleton class to give you an idea of how could be this new class.

```
# here goes import statement of the needed modules
import ABCD
from XYZ import xyz

# usually classes have CamelCase naming convention
class DiscountCurve:

    # the special __init__ method defines
    # how to construct instances of the class
    # so you need to identify the attributes you need to store
    # in the class defining a discount curve
    def __init__(self, ...):

        # then we want to add a method to compute the discount
        # factor at an arbitrary value date
        # using the data stored in the instance
    def df(self, param1, param2, ...):
        # the implementation can follow what we did in the
        # function we wrote last week but this time has to
        # use the class attributes
```

```
# finally we want a method to calculates the forward rate  
# based on the discount curve data stored in the instance  
def forward_rate(self, param1, param2, ...):  
    # here of course we can use the df method  
    # implemented above to calculate the forward rate
```



## Chapter 7

# Swaps and Bootstrapping

In this Chapter the Overnight Index Swap contract is reviewed and new class to represent it will be added to our financial module. Beside financial arguments another very important mathematical technique is introduced: the *bootstrapping*.

### 7.1 Payment Dates Generator

---

Before going to describe the Overnight Index Swap we need to develop a tool which helps us to generate list of dates (e.g. payment dates), a task that we will need to do often from now on. The function we are writing will go in `finmarkets` module and will be used by the classes describing various kind of contracts (this is essentially the function that was required in Ex. 3.5).

```
from datetime import date
from dateutil.relativedelta import relativedelta

def generate_swap_dates(start_date, n_months):
    dates = []
    for i in range(0, n_months, 12):
        dates.append(start_date + relativedelta(months=i))
    dates.append(start_date + relativedelta(months=n_months))

    return dates

print (generate_swap_dates(date.today(), 25))

[datetime.date(2020, 10, 20), datetime.date(2021, 10, 20), datetime.date(2022,
10, 20), datetime.date(2022, 11, 20)]
```

### 7.2 Overnight Index Swap

---

Interest rate swaps (IRS) are generally used to mitigate the risks of fluctuations of varying interest rates, or to benefit from lower rates.

Overnight Index Swaps (OIS) are a particular kind of IRS which pay a floating coupon, determined by overnight rate fixings over the reference periods, against a fixed coupon. By definition an OIS is defined by:

- a notional amount  $N$ ;
- a starting date  $d_0$ ;
- a sequence of payment dates  $d_1, \dots, d_n$ ;
- a fixed rate  $K$ .

For simplicity in the following we are assuming that the fixed and floating legs of our OIS have the same notional and payment dates, although this is not necessarily always the case in practice. We will always look at these products from the point of view of the **receiver of the floating leg**.

### 7.2.1 OIS Valuation

To evaluate the net present value (NPV) of such products the cash flows of each leg have to be calculated; today's NPV then is the sum of all the discounted cash flows.

#### Floating leg

At each payment date, the floating leg pays a cash flow determined as follows:

$$f_{\text{float}, i} = N \left\{ \prod_{d=d_{i-1}}^{d=d_i-1} \left( 1 + r_{\text{O/N}}(d) \cdot \frac{1}{360} \right) - 1 \right\}$$

Strictly speaking this formula is valid for an EONIA swaps (i.e.~for OIS swaps in EUR) other currencies might have different conventions. The  $\frac{1}{360}$  fraction appears because EONIA rates are quoted using the ACT/360 day-count convention. In addition we are making the simplifying assumption of ignoring weekends and holidays, so we assume that each overnight rate is valid for only one day. The sum of the discounted expected values of these cash flows is

$$\text{NPV}_{\text{float}} = \sum_{i=1}^n D(d_i) \mathbb{E}[f_{\text{float}, i}]$$

where  $D(d)$  is the discount factor with expiry  $d$ . On the other hand, by definition (see Section 6.3), the following relationship is also true

$$\mathbb{E}[f_{\text{float}, i}] = N \cdot \left( \frac{D_{\text{OIS}}(d_{i-1})}{D_{\text{OIS}}(d_i)} - 1 \right)$$

hence

$$\text{NPV}_{\text{float}} = N \cdot \sum_{i=1}^n D(d_i) \left( \frac{D_{\text{OIS}}(d_{i-1})}{D_{\text{OIS}}(d_i)} - 1 \right)$$

where  $D_{\text{OIS}}(d)$  is the discount factor implied by OIS prices (we will see how to derive it).

The correct curve to use for discounting the flows of a collateralized contract, like OIS, is the one associated with the collateral. Since OIS contracts are collateralized with cash, and cash accrues daily interest at the overnight rate, the OIS curve is itself the correct curve with which to discount the flows of an OIS contract ! So we have that  $D = D_{\text{OIS}}$  and the NPV simplifies to

$$\begin{aligned}
\text{NPV}_{\text{float}} &= N \cdot \sum_{i=1}^n [D(d_{i-1}) - D(d_i)] = \\
&= N \cdot [(D(d_0) - D(d_1)) + (D(d_1) - D(d_2)) + \dots + (D(d_{n-1}) - D(d_n))] \\
&= N \cdot [D(d_0) - D(d_n)]
\end{aligned} \tag{7.1}$$

### Fixed leg

The calculation for the fixed leg is simpler; each cash flow is equal to

$$f_{\text{fixed}, i} = N \cdot K \cdot \frac{d_i - d_{i-1}}{360}$$

so the NPV of the fixed leg is

$$\text{NPV}_{\text{fixed}} = N \cdot K \cdot \sum_{i=1}^n D(d_i) \frac{d_i - d_{i-1}}{360}$$

### 7.2.2 OvernightIndexSwap Class

Our ultimate goal is to take a series of Overnight Index Swap quotations, and determine the discount factors implied by their prices. To do this we will build a class to represent OIS and compute its value, given particular discount curve. Then we will use this class, put inside a numerical optimizer, to *invert* so that the implied discount factors can be determined from their prices (market quotes).

```

class OvernightIndexSwap:
    """
    OvernightIndexSwap: a class to valuate Overnight Index Swaps

    Attributes:
    -----
    notional: float
        Notional of the swap.
    payment_dates: list of datetime.date
        List of payment dates of the swap.
    fixed_rate: float
        Rate of the fixed leg of the swap.
    """
    def __init__(self, notional, payment_dates, fixed_rate):
        self.notional = notional
        self.payment_dates = payment_dates
        self.fixed_rate = fixed_rate

    def npv_floating_leg(self, discount_curve):
        """
        npv_floating_leg: computes the floating leg npv.

        Params:

```

```

        -----
        discount_curve: DiscountCurve
            Discount curve object used for npv calculation.
        """
        return self.notional * (discount_curve.df(self.payment_dates[0]) -
                                discount_curve.df(self.payment_dates[-1]))

def npv_fixed_leg(self, discount_curve):
    """
    npv_fixed_leg: computes the fixed leg npv.

    Params:
    -----
    discount_curve: DiscountCurve
        Discount curve object used for npv calculation.
    """
    npv = 0
    for i in range(1, len(self.payment_dates)):
        start_date = self.payment_dates[i-1]
        end_date = self.payment_dates[i]
        tau = (end_date - start_date).days / 360
        df = discount_curve.df(end_date)
        npv = npv + df * tau
    return self.notional * self.fixed_rate * npv
def npv(self, discount_curve):
    """
    npv: computes the total npv of the swap.

    Params:
    -----
    discount_curve: DiscountCurve
        Discount curve object used for npv calculation.
    """
    float_npv = self.npv_floating_leg(discount_curve)
    fixed_npv = self.npv_fixed_leg(discount_curve)
    return float_npv - fixed_npv

```

To test the newly developed class we need a discount curve. In the following example a fake curve will be defined, and then used with an OIS product.

```

from datetime import date
from finmarkets import DiscountCurve

ois = OvernightIndexSwap(
    # the notional, one million
    1e6,
    # the list of product dates,
    # i.e. the start date then the payment dates
    [date(2020, 1, 1), date(2020, 4, 1),

```



```

        date(2020, 7, 1), date(2020, 10, 1),
        date(2021, 1, 1)],
        # the fixed rate, 2.5%
        0.025)

# fake discount curve
curve = DiscountCurve(date(2020, 1, 1),
                      [date(2020, 1, 1), date(2021, 6, 1),
                       date(2022, 1, 1)],
                      [1.0, 0.98, 0.82])

ois.npv(curve)

105332.192377

```

## 7.3 Bootstrap Technique

As we said before we would like to determine a *real* discount curve starting from the market quotes of a set of Overnight Index Swaps with different maturities, this will be done via a technique called bootstrapping. This is the ABC of financial mathematics, since you almost always need a discount curve to price every contract. We are going to concentrate on EONIA swaps in order to build an EUR discount curve.

### 7.3.1 Building OIS instances

The first step involves getting data, the swap market quotes, and this is not actually as simple as it sounds.

The issue is that the EONIA swap market is over the counter (OTC) and it's not straightforward to access it. Unlike (some) listed futures, where anyone with a retail brokerage account can view and apply real time prices, to trade in the EONIA swap market you have to be a financial institution or at least a large company and have an agreement with a broker which operates in the market. One of the main brokers in the OIS market is ICAP, see Fig. 7.1. The underlying assumption is that market quotes represent the **fair price** of the OIS so they make the swap NPVs null (the fair price is an estimate of what a willing buyer would pay a willing seller for a given asset, assuming both have a reasonable knowledge of the asset's worth).

Though there exist some electronic platform in which market participants post bids and offers and other participants can apply them, in practice a lot of trading is still done over "voice", i.e. by phone or more commonly over chat. For convenience, however, Bloomberg provides a service which displays indicative real time rates as provided by a selection of relevant brokers. (*N.B. interest rate swap quotes vary from standard price quotes of commonly traded instruments, they can appear puzzling because the quotes are effectively interest rates*)

In the following we use a similarly created data-set (`ois_data.xlsx`) to derive our discount curve; with the help of the pandas module the data-set can be inspected:

```
import pandas, datetime
```

EONIA Rates up to 3YR				EONIA Rates 1-50YR				IMM FRA / EONIA SPREAD				ECB Dates EONIA				EUR Eonia vs USD OIS Basis Swap			
EONIA SWAPS																			
ICAP Global Menu -> ICAP EMEA -> Swaps -> OIS -> EUR -> EONIA Rates up to 3YR (GDCO 4963 10)																			
Term	Ask	Bid	Time	Term	Ask	Bid	Time												
1) 1 Week	-0.295	-0.395	07:00	16) 15 Month	-0.322	-0.372	11:46												
2) 2 Week	-0.297	-0.397	07:00	17) 18 Month	-0.319	-0.369	11:46												
3) 3 Week	-0.298	-0.398	07:00	18) 21 Month	-0.315	-0.365	11:46												
4) 1 Month	-0.325	-0.375	07:00	19) 2 Year	-0.309	-0.359	11:46												
5) 2 Month	-0.322	-0.372	07:00	20) 3 Year	-0.262	-0.312	11:46												
6) 3 Month	-0.323	-0.373	08:16	EONIA Forwards															
7) 4 Month	-0.324	-0.374	11:38	21) 1X2	-0.319	-0.369	07:00												
8) 5 Month	-0.324	-0.374	11:42	22) 2X3	-0.326	-0.376	11:45												
9) 6 Month	-0.324	-0.374	11:43	23) 1X4	-0.324	-0.374	11:38												
10) 7 Month	-0.324	-0.374	11:42	24) 2X5	-0.326	-0.376	11:43												
11) 8 Month	-0.323	-0.373	11:46	25) 3X6	-0.324	-0.374	11:46												
12) 9 Month	-0.323	-0.373	11:45	26) 6X12	-0.322	-0.372	11:46												
13) 10 Month	-0.323	-0.373	11:45																
14) 11 Month	-0.323	-0.373	11:46																
15) 12 Month	-0.322	-0.372	11:46																

Figure 7.1: Screenshot of market quotes from ICAP.

```
observation_date = datetime.date.today()

mq = pandas.read_excel('ois_data.xlsx')
mq.head()
```

```
   months  quote
0        1 -0.350
1        2 -0.347
2        3 -0.348
3        4 -0.350
4        5 -0.350
```

Next we could convert the data-set into a dictionary for later usage or use the DataFrame directly, it is just matter of taste.

Let's say we want to build a 15 months swap instance using data contained in ois\_data file. Be careful when doing this operation and double check the units of rates, quotes, etc... in this case for example quotes are expressed in percent so you need to multiply them by 0.01 before using them. Another detail to check is that 15 months quote is not the fifteenth entry in the DataFrame (actually it is the twelfth).

```
ois = OvernightIndexSwap(1e6,
                          [date(2019, 10, 23),
                           date(2020, 10, 23),
                           date(2020, 1, 23)],
                          mq['quote'].tolist()[12]*0.01)

# print the last payment date
# (15 months after obs date)
ois.payment_dates[-1]

datetime.date(2020, 1, 23)
```

Clearly to use the npv method to calculate the OIS' NPV we need a discount curve with which to evaluate it and here comes to hand the bootstrapping technique !

### 7.3.2 Constructing the Yield Curve

Keep aside for a moment our swaps and introduce the *bootstrap algorithm*. In finance, bootstrap is a method for constructing a (zero-coupon) fixed-income yield curve from the prices of a set of coupon-bearing products, e.g. bonds and swaps. The term structure of spot returns is obtained from the bond yields by solving for them recursively, by forward substitution: this iterative process is what is called the bootstrap method. The usefulness of bootstrap is that using only a few carefully selected zero-coupon products, it becomes possible to derive swap forward and spot rates for all maturities given the solved curve.

To illustrate bootstrapping let's consider the following example which can be partially solved analytically: we have some coupon paying bond (coupon of 4%, 5%, 6%, 7% and 8% respectively) with maturities ranging from 1 to 5 years, each having a value of €100 and traded at par. To determine the zero-coupon yield curve proceed as follows:

1. at the end of the first year the 1<sup>st</sup> bond will pay a coupon of €4 (= €100 \* 4%) plus the principal amount (= €100) which sums up to €104 while the bond is trading at €100. Therefore, the implied 1-year spot *fair* rate  $S_{1y}$  can be calculated as,  $€100 = €104 / (1 + S_{1y})$ ;
2. at the end of second year the sum of the cash flows of the 2<sup>nd</sup> bond can be compared to its trading price to compute the 2-year spot rate  $S_{2y}$  as  $€100 = €5 / (1 + S_{1y}) + €105 / (1 + S_{2y})^2$ , using the previously derived value of  $S_{1y}$ ;
3. at the end of third year the sum of the cash flows of the 3<sup>rd</sup> bond can be compared to its trading price to calculate the 3-year spot rate  $S_{3y}$  as  $€100 = €6 / (1 + S_{1y}) + €6 / (1 + S_{2y})^2 + €106 / (1 + S_{3y})^3$ , using  $S_{1y}$  and  $S_{2y}$  computed before;
4. repeat the same reasoning for the other bonds.

Putting all together we can construct a system of equations (now omitting the currency symbol for simplicity):

$$\left\{ \begin{array}{l} 100 = \frac{104}{(1 + S_{1y})} \\ 100 = \frac{5}{(1 + S_{1y})} + \frac{105}{(1 + S_{2y})^2} \\ 100 = \frac{6}{(1 + S_{1y})} + \frac{6}{(1 + S_{2y})^2} + \frac{106}{(1 + S_{3y})^3} \\ 100 = \frac{7}{(1 + S_{1y})} + \frac{7}{(1 + S_{2y})^2} + \frac{7}{(1 + S_{3y})^3} + \frac{107}{(1 + S_{4y})^4} \\ 100 = \frac{8}{(1 + S_{1y})} + \frac{8}{(1 + S_{2y})^2} + \frac{8}{(1 + S_{3y})^3} + \frac{8}{(1 + S_{4y})^4} + \frac{108}{(1 + S_{5y})^5} \end{array} \right. \quad (7.2)$$

This system can be solved quite easily: from the first equation can be derived  $S_{1y}$ , from the second  $S_{2y}$ , from the third  $S_{3y}$  and so on. So

$$100 = 104 / (1 + S_{1y}) \quad \Rightarrow \quad S_{1y} = 104 / 100 - 1 = 4\%$$

Moving to the second equation:

$$100 = 5/(1 + 0.04) + 105/(1 + S_{2y})^2 \Rightarrow S_{2y}^2 + 2S_{2y} - 0.103030 = 0$$

$$S_{2y} = -1 \pm \sqrt{1 + 0.103030} = \begin{cases} -2.05023 \\ 0.0503 \end{cases}$$

where the first solution has been discarded because negative.

From the third one on it is not as simple to solve them analytically since involve third order (or more) equations, anyway it is possible to solve them numerically. Assume we have found all the rates up to the fourth year (they are reported in Table 7.1) and let's try to determine the last one.

years	coupon rate	bond price	spot rate
1	1.00 %	€100	4.00%
2	2.00 %	€100	5.03%
3	3.00 %	€100	6.08%
4	4.00 %	€100	7.19%
5	5.00 %	€100	???

Table 7.1: Table reporting maturity, coupon, bond price and implied spot rate for the example outlined in the text.

The last column of Table 7.1 provides with the terms to fill the zero-coupon yield curve. To solve the last equation numerically we can use the `scipy.optimize.brentq` function which finds the zeros of a user-defined function given a validity interval. In Figure 7.2 the function to determine the 5 year rate expressed in the last of Eqs. 7.2 is shown. From the plot we expect the rate to be around 8%.

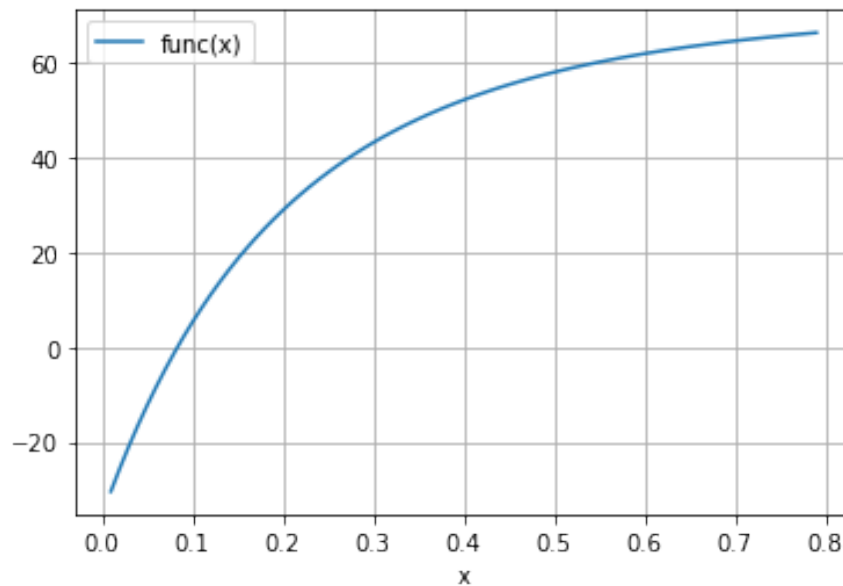


Figure 7.2: Plot of the discounted cash flow of bond 5 as a function of the 5 year spot rate.

```

from scipy.optimize import brentq

def func(x):
    return 100 - 8/(1+0.04) - 8/(1+0.0503)**2 - 8/(1+0.0608)**3
        - 8/(1+0.0719)**4 - 108/(1+x)**5

a = brentq(func, 0, 0.10)
print ("5y rate: {:.4f}".format(a))

5y rate: 0.0836

```

The very same mechanism can be generalized and extended to more maturities to get a more detailed yield curve. In general terms the previous system can be written as:

$$\begin{cases} f_1(S_1, p_1) = 0 \\ f_2(S_1, S_2, p_2) = 0 \\ f_3(S_1, S_2, S_3, p_3) = 0 \\ f_4(S_1, S_2, S_3, S_4, p_4) = 0 \\ \dots \end{cases}$$

where  $S_i$  are the unknown spot rates and  $p_i$  the prices of the considered products. The iterative procedure we have applied before exploits the first equation to find  $S_1 = f_1^{-1}(p_1)$ , the second to find  $S_2 = f_2^{-1}(S_1, p_2)$  and so on and so forth; this algorithm works since each equation will determine exactly one *free* spot rate which is not already determined by the others.

### 7.3.3 Bootstrap as Minimization Problem

We can now describe the bootstrapping algorithm in general terms as follows:

1. define a set of yielding products , these will generally be coupon-bearing bonds;
2. derive discount factors for the corresponding terms;
3. *bootstrap* the zero-coupon curve, successively calibrating the curve such that it returns the prices of these inputs.

Instead of iteratively finding the solution of each equation as before, we could define a vector of spot rates  $\mathbf{S} = (S_1, S_2, S_3, \dots)$  seeking for a particular  $\hat{\mathbf{S}}$  which solves the following equation:

$$F = f_1^2(\hat{S}_1) + f_2^2(\hat{S}_1, \hat{S}_2) + f_3^2(\hat{S}_1, \hat{S}_2, \hat{S}_3) + f_4^2(\hat{S}_1, \hat{S}_2, \hat{S}_3, \hat{S}_4) + \dots = 0$$

Under this terms bootstrapping can be considered as a minimization algorithm, indeed we need to find  $\hat{\mathbf{S}}$  which *minimize*  $F$ , or makes it as close as possible to 0. Notice how each  $f_i$  is squared since we want all of them to be minimized and not only  $F$  globally (without the squared there may be cancellation effects between the terms of the sum).

### 7.3.4 Minimization Algorithm

A minimization algorithm follows these steps:

- define an *objective function* i.e. the function that is actually minimized to reach our goal;

- set the initial value of the unknown parameters and their range of variability;
- the minimizer will compute the objective function value;
- then it will move the parameter values in such a way to find a smaller value of the objective function (e.g. following the derivative w.r.t. each parameter);
- if constraints are defined, they will be considered in the previous step;
- the last three steps will be repeated until further variations of the  $x$  values won't change significantly the objective function (i.e. we have found a minimum of the function so the minimization process is completed !).

Let's see with a couple of example how minimization can be implemented in python using the function `scipy.optimize.minimize`.

### Example

Find the dimensions that will minimize the cost to manufacture a circular cylindrical can of volume,  $330 \text{ cm}^3$ , see Figure 7.3.

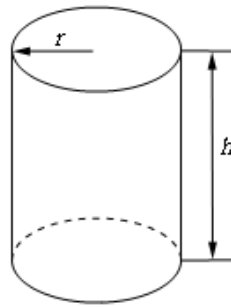


Figure 7.3: Graphical representation of the *can* minimization example.

Clearly to minimize the costs the company needs to reduce the can surface, given the required volume.

$$S = 2\pi rh + 2 \cdot (\pi r^2)$$

On the other hand we want the volume to be  $330 \text{ cm}^3$  so we can remove  $h$  from the previous equation:

$$V = \pi r^2 h = 330 \quad \implies \quad h = \frac{330}{\pi r^2}$$

So in the end the surface function to be minimized is:

$$S = 2\pi rh + 2 \cdot (\pi r^2) = \frac{2 \cdot 330}{r} + 2 \cdot (\pi r^2)$$

So we implement the objective function,  $x[0]$  is the can radius:

```
from math import pi

def of(x):
    return 2*330/x[0] + 2*pi*x[0]**2
```

Set the limits to our unknown variable and its initial value:

```
x0 = [1]
bounds = [(0.01, 100)]
```

Finally we run the minimization:

```
r = minimize(of, x0, bounds=bounds)
print (r)

    fun: 264.356810914805
 hess_inv: <1x1 LbfgsInvHessProduct with dtype=float64>
   jac: array([5.68434189e-06])
message: b'CONVERGENCE: NORM_OF_PROJECTED_GRADIENT_<=_PGTOL'
  nfev: 24
   nit: 9
status: 0
success: True
     x: array([3.7449385])
```

So to minimize the cost the company should produce cans with a radius of about 3.745 cm (I suspect that Coke have done a similar calculation...).

### Example with Constraint

We are going to fence in a rectangular field. If we look at the field from above the cost of the vertical sides are \$10/m, the cost of the bottom is \$2/m and the cost of the top is \$7/m. If we have \$700 determine the dimensions of the field that will maximize the enclosed area, see Fig. 7.4.

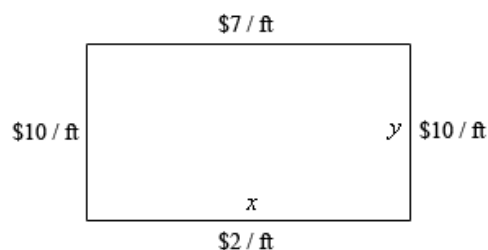


Figure 7.4: Graphical representation of the *field* minimization example.

In this example there are two differences with respect to the previous one:

- we want to maximize a quantity (not minimize);
- there is a constraint (we have a limited amount of money).

So let's repeat the steps as before. The objective is to maximize the enclosed area  $A$  but we are able to just minimize so we can define in the objective function the quantity  $-A$ , if we minimize it we will maximize the area  $A$ . Define the length and the width of the field with  $x[0]$  and  $x[1]$  (items of the list  $x$ ):

```
def of(x):
    return - x[0]*x[1]
```

Now we can set the boundaries for length and width and their initial values (1 m each):

```
x0 = [1, 1]
bounds = [(0.01, 100) for _ in range(len(x0))]
```

We have also to impose the constraint on the money. This is done by defining a function that compute the money spent with the fence and compare it to \$700. The constraint is passed to the minimizer with a dictionary which has two keys: `type` with value `eq` (like equality) since we want to spend all of our available money so the fence has to cost \$700

$$\text{fence cost} = l \cdot 10 + l \cdot 10 + w \cdot 2 + w \cdot 7 = 700$$

$$700 - l \cdot 10 - l \cdot 10 - w \cdot 2 - w \cdot 7 = 0$$

, fun whose value is the constraint function.

```
def cons(x):
    return 700 - x[0]*20 - x[1]*2 - x[1]*7

constraints = {'type': 'eq', 'fun': cons}
```

Now we can call the minimizer.

```
r = minimize(of, x0, bounds=bounds, constraints=constraints)
print (r)

fun: -680.5555555555482
jac: array([-38.88889313, -17.5      ])
message: 'Optimization terminated successfully.'
nfev: 16
nit: 4
njev: 4
status: 0
success: True
x: array([17.49999818, 38.88889293])
```

So the field will come out 17.5m long and 38.9m wide.

### 7.3.5 OIS Example

Back to our Overnight Index Swap, the general idea here is to get the discount curve  $\mathcal{C}$  such that it prices correctly each OIS by minimizing the sum of their NPV squared (our  $f_i$ ):

$$\min_{\mathcal{C}} \left\{ \sum_{i=1}^n \text{NPV}(\text{OIS}_i, \mathcal{C})^2 \right\}$$

A discount curve is characterized by pillar dates and the corresponding discount factors. The description of the problem we have given above does not, in theory, specifies any constraint on the number of pillar dates of the discount curve  $\mathcal{C}$ . However, the pillar dates determine the number of unknown variables (i.e. the dimensionality  $N$  of the optimization problem). A curve with  $N$  pillar dates has  $N$  discount factors (note that the first discount factor with value date equal to the today date, is constrained to 1). **In practice, therefore, it makes sense to choose the pillar dates in such a way that there are exactly the right number of degrees of freedom in the optimization**



**to match data.** So the natural choice is to choose the pillar dates of the discount curve equal to the set of expiry dates of the swaps.

Once we've fixed  $\mathbf{d}$  to be a vector of pillar dates equal to the expiry dates of the swaps, and we use the notation  $\mathbf{x}$  to represent the vector of pillar discount factors, then the problem becomes:

$$\min_{\mathbf{x}} \left\{ \sum_{i=1}^N \text{NPV}(\text{OIS}_i, \mathcal{C}(\mathbf{d}, \mathbf{x}))^2 \right\}$$

which is our optimization problem (**to find the minimum of the above expression as a function of  $\mathbf{x}$** ) that can be solved using one of the available numerical optimization routines in python.

So first let's create the swaps according to all the available market quotes and also the pillar dates of our final discount curve:

```
from finmarkets import generate_swap_dates

observation_date = date(2019, 10, 23)
pillar_dates = [observation_date]
swaps = [] # container of the OIS objects

for i in range(len(df)):
    swap = OvernightIndexSwap(1e6,
                              generate_swap_dates(
                                  observation_date,
                                  mq['months'].tolist()[i],
                                  0.01 * mq['quote'].tolist()[i])

    swaps.append(swap)
    pillar_dates.append(swap.payment_dates[-1])

# this shouldn't be necessary if the original
# list of market quotes is sorted
pillar_dates = sorted(pillar_dates)
```

So let's implement the method with the swaps we have just created, of course we don't need to write our minimisation algorithm since we can use the one provided by python which is defined in `scipy.optimize`, function `minimize`.

- define the objective function: the sum of the squared NPVs of the OIS

```
def objective_function(x):
    curve = DiscountCurve(observation_date,
                          pillar_dates,
                          x)

    sum_sq = 0.0
    for swap in swaps:
        sum_sq += swap.npv(curve) ** 2
    return sum_sq
```

- set the initial value of the discount factors ( $x_i^0$ ) to 1 with a range of variability  $[0.01, 10]$ , in addition the first element of the list, today's discount factor, will be fixed to 1 (variability  $[1, 1]$ )

```
x0 = [1.0 for i in range(len(pillar_dates))]  
  
bounds = [(0.01, 10.0) for i in range(len(pillar_dates))]  
bounds[0] = (1.0, 1.0)
```

- finally we can launch the minimizer to find the discount factors ( $x$ )

```
from scipy.optimize import minimize  
  
result = minimize(objective_function, x0, bounds=bounds)  
print (result)
```

```
fun: 0.000819919032900304  
hess_inv: <34x34 LbfgsInvHessProduct with dtype=float64>  
jac: array([ 6.58948735e+05, -1.58720803e+01, -6.53143264e+01,  
-1.03323232e+02,  
-1.26050260e+02, -1.31748898e+02, -1.20374599e+02, -9.15399651e+01,  
-4.24363322e+01,  2.44903182e+01,  1.14345243e+02,  2.22002243e+02,  
-3.72021700e+00,  4.21398633e+01,  4.21787852e+01,  4.22369487e+01,  
 4.23327026e+01,  4.31814758e+01,  4.44924460e+01,  4.62078978e+01,  
 4.82906823e+01, -3.69972738e+00, -1.42454702e+00,  7.53771932e-01,  
 2.79741018e+00,  4.62896699e+00,  6.24844054e+00,  9.93101553e+00,  
 1.31122434e+01,  1.42880909e+01,  1.48279215e+01,  1.50787019e+01,  
 1.43267935e+01,  1.38451324e+01])  
message: b'CONVERGENCE: REL_REDUCTION_OF_F_<=_FACTR*EPSMCH'  
nfev: 840  
nit: 7  
status: 0  
success: True  
x: array([1.          , 1.00030147, 1.00058831, 1.00089012, 1.00119726,  
1.00147996, 1.00178743, 1.00208107, 1.00238467, 1.00267865,  
1.00298261, 1.00327737, 1.00357104, 1.00357104, 1.00355063,  
1.00352002, 1.00346901, 1.00302007, 1.00232627, 1.00141821,  
1.00031629, 0.99911234, 0.99790839, 0.99675545, 0.99567393,  
0.99470465, 0.9938476 , 0.99189884, 0.99021534, 0.98959296,  
0.98930728, 0.98917464, 0.98957256, 0.98982763])
```

Printing the result gives us the a lot of information about the minimisation just performed, the most useful are:

- func: the value of the objective function at the last iteration;
- message: the summary message from the algorithm (if it is CONVERGENCE is OK);
- success: the name is self explanatory;
- x: the vector of unknown parameters that have been optimised.

Another useful check to perform in order to understand if everything went fine, is the comparison of the objective function with the initial guessed parameters and at the end of the minimization.

```
print ("Initial objective function value ", objective_function(x0))
print ("Final objective function value ", objective_function(result.x))
```

```
Initial objective function value  931188216.6666666
Final objective function value  0.000819919032900304
```

The objective function at the end of the minimisation is not exactly 0 (and rarely it will be) but its value is small enough for us to be satisfied, we started with  $10^{10}$  and now it is  $10^{-4}$  so 14 orders of magnitude smaller. This means that with the derived discount curve the NPV's of our OIS won't be identically 0 but so small that we can consider them as they were.

It can be very useful to also look at some diagnostic plots to check if the minimization was successful. Figure 7.5 reports on the left the objective function value as a function of discount factor ( $x_1$ ); clearly we have found a minimum (the orange point represent  $x_1$  value at the end of the minimization). On the right the value of the objective function at each iteration is shown instead, its value is decreasing dramatically (notice that the y axis is drawn in log scale).

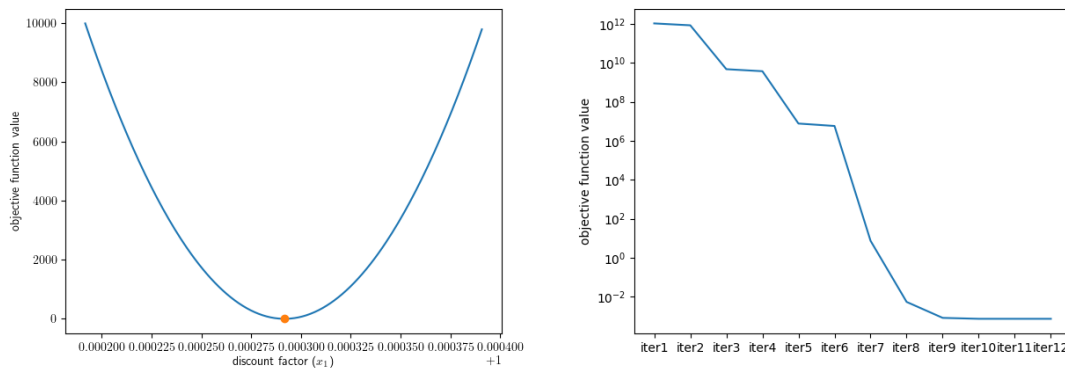


Figure 7.5: Diagnostic plots for the minimization algorithm. On the left the objective function value as a function of the discount factor  $x_1$ , on the right the objective function value as a function of the iteration number (the orange point represent  $x_1$  value at the end of the minimization).

Finally we can create the discount curve implied by the market quote of our swaps (see Fig. 7.6) and try to compute some implied rate.

```
from math import log
curve = DiscountCurve(observation_date, pillar_dates, result.x)

d = date(2059, 11, 23)
print ("40y df: {}".format(curve.df(d)))
print ("40y rate: {}".format(-log(curve.df(d)) / 40))

40y df: 0.9891780176191146
40y rate: 0.0002720241491103593
```

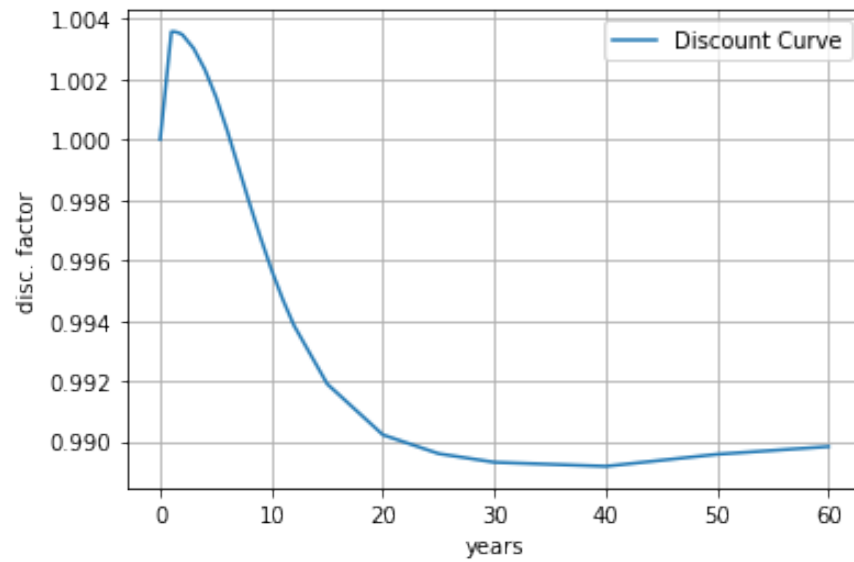


Figure 7.6: Plot of the discount curve implied by Overnight Index Swap market quotes.