

# INTRODUZIONE AGLI SCRIPT PYTHON PER L'ANALISI DATI

## SOMMARIO

L'analisi dati richiede spesso di eseguire ripetutamente operazioni semplici, ad esempio nel caso di calcolo di medie o deviazioni standard, oppure nel caso di fit numerici. È in questi casi che il calcolatore ci può aiutare a velocizzare le operazioni. Un semplice *script*, basato su un linguaggio ad alto livello, è quello che serve nella quasi totalità dei casi.

Nell'esperienza di oggi vedremo alcuni esempi di base che possono essere utili nelle future esperienze di laboratorio. Utilizzeremo il linguaggio Python, noto per essere semplice da usare, dotato di un'ampia libreria standard e di librerie aggiuntive (scipy, numpy, matplotlib, etc.) che lo rendono molto pratico per l'analisi dati.

## L'AMBIENTE DI SVILUPPO

Python può essere usato sia in maniera interattiva, aprendo il suo interprete e inserendo le istruzioni una alla volta, che per eseguire script salvati su file di testo. Il primo caso si presta ad operazioni immediate da fare una volta sola. Il secondo, su cui ci soffermeremo, è più adatto quando lo script è più complesso di qualche semplice riga e quando si vuole avere la libertà di fare piccole modifiche allo script ed eseguirlo più volte senza troppi sforzi.

Potete usare un qualunque editor di testo (emacs, gedit, vi, etc.) per scrivere il vostro codice, ed eseguirlo nel terminale a riga di comando, con comandi come “python myscript.py”. Esistono anche ambienti di sviluppo integrati come Pyzo (disponibile nei PC in laboratorio) che forniscono sia dei tab per scrivere e salvare i vostri script, gli strumenti per eseguire il codice, un insieme di librerie pre-installate ed una shell python interattiva.

## IL TERMINALE

È consigliato crearsi *directory*, (ovvero una cartella) dedicata per ciascun gruppo ed esperienza in cui lavorare e mettere i propri file (script, dati, grafici, etc.). Partendo dai comandi base di terminale del sistema operativo Linux:

- `pwd` restituisce sullo schermo il nome della directory corrente;
- `ls` elenca i file contenuti in una certa directory;
- `cd` cambia directory (es. “`cd plasduino`”);
- `mkdir` crea directory (es. “`mkdir esperienza_dadi`”);
- `cp` copia file (es. “`cp plasduino/data/dati.txt mydir`”);
- `mv` sposta o rinomina un file (es. “`mv dati.txt dati_1.txt`”);
- `rm` cancella file (es. “`rm dati.txt`”)

Apriamo un terminale, controlliamo in che cartella siamo (con `pwd`), creiamo la nostra cartella di lavoro (con `mkdir`) e entriamo (con `cd`). A questo punto facciamo partire Pyzo e creiamo un nuovo script nella nostra cartella di lavoro.

## PRIMO PROGRAMMA

È tradizione, iniziare ad apprendere i fondamenti della sintassi di un linguaggio di programmazione, con un semplice programma che restituisca una stringa di testo. Nel caso di Python, si scrive un output su terminale semplicemente con la funzione `print`<sup>1</sup>:

```
1 # primo programmino
2 print('Hello world!')
```

## GIUSTO UN PO' PIÙ COMPLESSO

Notare che `print` restituisce sempre una stringa di testo e quindi prova a convertire in questo formato anche variabili che testo non sono. Ecco un esempio con definizione e manipolazione di variabili, in cui si vede anche come importare la libreria delle funzioni matematiche:

```
1 # definisco x e y come interi
2 x = 2
3 y = 3
4 # stampo la somma
5 print (x+y)
6 #importo una libreria con funzioni matematiche
7 # ne vedremo dopo le caratteristiche
8 import numpy
9 # calcolo la radice quadrata della somma
10 z = numpy.sqrt(x+y)
11 # stampo la radice quadrata della somma
12 print (z)
13 # forzando solo 3 cifre dopo la virgola
14 print ('%.3f' % z)
15 # forzando ad intero
16 print (int(z))
```

Cosa succede se c'è un errore nel codice? Ad esempio se proviamo a fare la radice quadrata di una lettera:

```
1 import numpy
2 x = numpy.sqrt('a')
```

Otteniamo un messaggio di errore con indicata la linea dove è stato trovato l'errore ed il suo tipo. In questo caso “`TypeError: Not implemented for this type`”, ovvero la funzione `sqrt` non funziona con le stringhe di testo, ma si aspetta numeri. È bene leggere con cura tali messaggi (*Traceback*) di solito sono abbastanza chiari da indicare bene quale errore abbiamo fatto.

1. In tutte le lingue del mondo: <http://helloworldcollection.de/>

## LEGGERE E USARE UN FILE DI DATI

Può capitare che i dati da analizzare siano scritti in un file di testo, vediamo quindi come leggerlo. Iniziamo dal crearne uno con un po' di dati, ad esempio due colonne separate da una tabulazione (tasto "tab")

```
1 #x      y
2 1.0     2.08
3 2.0     4.44
4 3.0     6.09
5 4.0     6.88
```

Proviamo a leggerlo e a riempire due liste con numeri *float* (ovvero in virgola mobile). È fondamentale l'indentazione, ovvero che righe di comandi in uno stesso gruppo siano allineate verticalmente.

```
1 # apriamo il file
2 fin = open('dati.txt')
3 # cosa succede se faccio print?
4 print(fin)
5 # andiamo attraverso ogni linea e
6 # riempiamo una lista di numeri:
7 ListaX = []
8 ListaY = []
9 for line in fin:
10     if not line.startswith('#'):
11         [x,y] = line.split('\t') # separiamo sui tab
12         print(x,y) # diamo un'occhiata al contenuto
13         ListaX.append(float(x)) # forziamo a float
14         ListaY.append(float(y))
15
16 # vediamo il risultato:
17 print(ListaX)
18 print(ListaY)
```

## MEDIA ARITMETICA

Vogliamo definire una funzione per il calcolo della media di una lista di dati.

$$m = \frac{1}{n} \sum_{i=1}^n x_i \quad (1)$$

```
1 # definiamo la funzione
2 def media(x):
3     n = len(x) # lunghezza della lista
4     media = 0. # inizializziamo la media
5     for i in x:
6         media += i
7     media /= n
8     # restituiamo la media
9     return media
10
11 # inventiamoci un set di valori
12 dati = [1,2,3,4,5,6]
13 # calcoliamone la media
14 print(media(dati))
```

Notare che la media è inizializzata con un float, cosa succede se invece uso un numero intero (*int*)?

## DEVIAZIONE STANDARD

Proviamo ora ad aggiungere alla media, il calcolo della deviazione standard:

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - m)^2} \quad (2)$$

```
1 # ci serve la radice quadrata
2 from numpy import sqrt
3
4 # definiamo la funzione
5 def getstats(x):
6     # media come prima
7     n = len(x) # lunghezza della lista
8     media = 0. # inizializziamo la media
9     for i in x:
10         media += i
11     media /= n
12     # ora somma dei quadrati:
13     ss = 0.
14     for i in x:
15         ss += (i-media)**2
16     # normalizzazione:
17     stdev = sqrt(ss/(n-1))
18     # in output i 2 valori
19     return media, stdev
20
21 # inventiamoci un set di valori
22 dati = [1,2,3,4,5,6]
23 # e usiamo la funzione
24 print(getstats(dati))
```

Notiamo che nella funzione precedente abbiamo fatto 2 *loop*, ovvero due cicli sullo stesso set di dati. C'è un modo per rendere la funzione più veloce evitando il secondo *loop*? Possiamo provare a manipolare la formula della deviazione standard:

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i^2 + m^2 - 2x_i m)} \quad (3)$$

$$= \sqrt{\frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 + n m^2 - 2m \sum_{i=1}^n x_i \right)} \quad (4)$$

$$= \sqrt{\frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - n * m^2 \right)} \quad (5)$$

e notiamo che basta la somma dei valori e dei suoi quadrati per costruire media e deviazione standard. Scriviamo quindi una nuova funzione:

```
1 # ci serve la radice quadrata
2 from numpy import sqrt
3
4 # definiamo la funzione
5 def getstats2(x):
6     # media come prima
7     n = len(x) # lunghezza della lista
8     s = 0. # inizializziamo somma
9     ss = 0. # inizializziamo somma dei quadrati
```

```

10 for i in x:
11     s += i
12     ss += i*i
13 # calcoliamo media:
14 media = s/n
15 # calcoliamo stdev:
16 stdev = sqrt((ss - n*media**2)/(n-1))
17 # in output i 2 valori
18 return media, stdev
19
20 # inventiamoci un set di valori
21 dati = [1,2,3,4,5,6]
22 # e usiamo la funzione
23 print(getstats2(dati))

```

Notate differenze tra il risultato delle due funzioni? Sapendo che la seconda è più veloce, potete immaginare casi in cui è preferibile la prima?

## NUMPY

Una libreria molto usata per il calcolo numerico e scientifico in Python è NumPy (<http://www.numpy.org/>). Il suo scopo è principalmente la manipolazione di *array* multidimensionali. Un array (unidimensionale) è una sequenza ordinata di numeri, quindi abbastanza simile ad una lista, ma con differenze fondamentali. Vediamo, nell'esempio seguente, alcune caratteristiche degli array di numpy:

```

1 # cominciamo importando numpy
2 import numpy as np
3 # creiamo quindi il nostro primo array:
4 a = np.array([2,3,4])
5 # e proviamo a stamparlo
6 print(a)
7 # Prima cosa importante: il tipo di variabile
8 # e' fissato per tutti per tutti gli elementi
9 print(a.dtype) # in questo caso un int
10 # tipi diversi si possono specificare:
11 b = np.array([2,3,4], 'd')
12 print(b.dtype)
13
14 # il metodo append() esiste solo per le liste:
15 l = [2,3,4]
16 l.append(5)
17 print(l)
18 a.append(5) # Riuscite a capire il Traceback?
19 # per aggiungere elementi la sintassi e':
20 a = np.insert(a, len(a), 5.2)
21 print(a) # Perché compare 5 e non 5.2?
22
23 # le operazioni tra array sono molto immediate:
24 print(a+a)
25 print(a**2)
26 # ma attenzione a usare il metodo giusto:
27 A = np.array([[1,1],
28               [0,1]])
29 B = np.array([[2,0],
30               [3,4]])
31 print(A*B) # elementwise product
32 print(A.dot(B)) # matrix product

```

## MEDIA E DEVIATION STANDARD CON NUMPY

Come prima, ma stavolta usando `array`:

```

1 import numpy as np
2
3 def getstatsNP(sample):
4     mean = sample.sum()
5     variance = (sample**2).sum()
6     n = len(sample)
7     mean /= n
8     variance = (variance - n*mean**2.0)/(n - 1)
9     return mean, np.sqrt(variance)
10
11 # Vediamo i risultati:
12 dati = np.array([1,2,3,4,5,6]) # con int
13 print(getstatsNP(dati))
14 dati = np.array([1,2,3,4,5,6], 'd') # con float
15 print(getstatsNP(dati))

```

## GRAFICI CON MATPLOTLIB

Esistono diverse librerie grafiche che consentono la creazione di grafici di vario tipo. In questa esperienza vedremo brevemente alcuni esempi con la libreria `matplotlib`. Cominciamo con un semplice grafico cartesiano:

```

1 # cominciamo ad importare la libreria
2 import numpy as np
3 import matplotlib.pyplot as plt
4 # inventiamoci un set di dati
5 x = np.arange(1,9,1)
6 y = np.array([2.08, 4.44, 6.09, 6.88, 8.61,
7               9.79, 11.9, 12.6])
8 xe = np.array(8*[0.1])
9 ye = np.array(8*[0.3])
10 # creiamo uno scatter plot
11 plt.errorbar(x, y, ye)
12 # magari mettiamo anche gli errori in x
13 plt.errorbar(x, y, ye, xe, 'o')
14 # e mostriamo il grafico finale
15 plt.show()

```

È bene imparare a formattare i grafici in modo che siano meglio comprensibili. Proviamo a trarre ispirazione dalle codice seguente, potete intuire a cosa servono le varie righe?

```

1 plt.rc('font', size = 18)
2 plt.title('Law of motion', y = 1.02)
3 plt.xlabel('t [s]')
4 plt.ylabel('s [m]', labelpad = 25)
5 plt.xlim(0, 10)
6 plt.ylim(0, 18)
7 plt.grid(color = 'gray')

```

Infine sarebbe tutto inutile se non foste in grado di salvare il grafico finale:

```

1 plt.savefig('mioplot.pdf') # o .png etc.

```

Oltre agli *scatterplot* ci sono altri tipi di grafici disponibili. Un altro esempio è il grafico a barre:

2. In realtà NumPy offre anche funzioni predefinite per il calcolo delle statistiche campione

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 # inventiamo i x e y (valori e occorrenze)
4 valori      = np.array([1,2,3,4,5])
5 occorrenze  = np.array([.2,.5,2,.6,.7])
6
7 plt.title('Law of motion', y = 1.02)
8 plt.xlabel('Valori')
9 plt.ylabel('Occorrenze')
10 plt.xlim(0, 6)
11 plt.ylim(0, 3)
12 plt.grid(color = 'gray')
13 plt.bar(valori, occorrenze, width = 1)
14 plt.savefig('barplot.png')
```

Notare come sono “centrate” le barre sull’asse x