

## Relazione per il corso di Data Science

Liam Cavini  
Semestre Invernale 2024/2025

2° Foglio, Ottimizzazione  
23/10/2024

## Esercizio 1 - Ottimizzazione non vincolata

(a) Per trovare i minimi della funzione:

$$f(x, y) = x^2 - 2x + x^2y^2 - 2xy$$

si calcola il gradiente e si trovano i punti critici. Si studia successivamente se la hessiana valutata in questi punti è definita o semidefinita positiva. Si trova così che i punti critici sono  $(1, 1)$  e  $(0, -1)$ , e che il primo di questi corrisponde ad un minimo locale ed il secondo ad un punto di sella.

Partendo da  $(0, 0)$  il metodo della discesa del gradiente (con tasso di apprendimento  $\alpha = 0.2$ ) trova il minimo senza problemi, come si può osservare in Figura 1.

Il metodo di Newton invece in questo caso fallisce, in un solo step infatti raggiunge il punto di sella, come mostrato sempre in Figura 1. Se si calcola la hessiana esplicitamente nel punto  $(0, 0)$  si trova che:

$$H(0, 0) = \begin{bmatrix} 0 & -0.5 \\ -0.5 & -0.5 \end{bmatrix}$$

Moltiplicando la hessiana per il gradiente in  $(0, 0)$  (il quale risulta essere  $(-2, 0)$ ), si osserva infatti che il metodo di Newton non aggiorna il parametro  $x$  nella prima iterazione, aggiornando invece  $y$  ad un valore più piccolo, allontanandosi di conseguenza minimo.

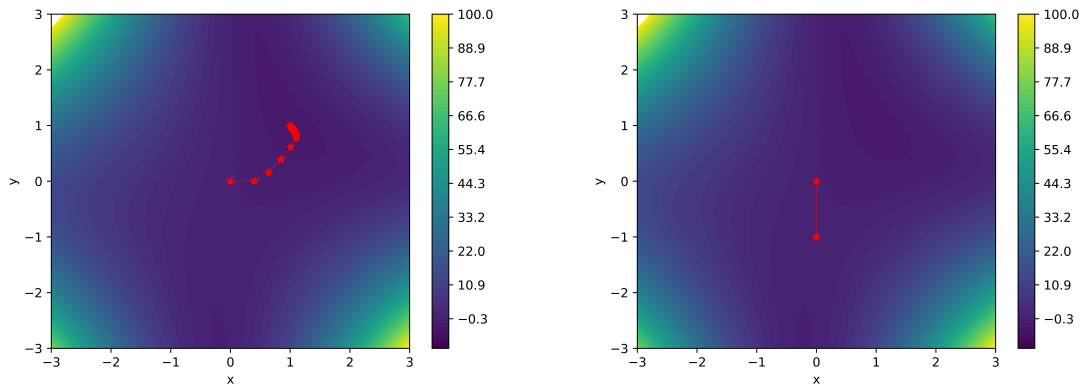


Figura 1: Le immagini mostrano una planimetria della funzione. I tratti rossi indicano il percorso compiuto dalla discesa del gradiente che inizia dal punto  $(0, 0)$ . A sinistra abbiamo il risultato ottenuto con il metodo classico, con tasso di apprendimento pari a 0.2, che raggiunge il minimo globale. A destra il risultato ottenuto col metodo di Newton, che si ferma sul punto di sella.

Il codice usato per la discesa del gradiente è il seguente:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def gradient(pos):
5     x = pos[0]
6     y = pos[1]
7     return np.array([2*x-2+2*x*y**2-2*y, 2*x**2*y -2*x])
8
9 def update(pos, learning_rate):
```

```

10     pos -= learning_rate * gradient(pos)
11
12     iteration_number = 10000
13     initial_pos = np.array([0,0])
14     positions = np.zeros((iteration_number+1,2))
15     positions[0] = initial_pos
16
17     learning_rate = 0.2
18
19     for i in range(iteration_number):
20         positions[i+1] = positions[i]
21         update(positions[i+1], learning_rate)
22
23     #matplotliblib
24     X, Y = np.meshgrid(np.linspace(-3, 3, 128), np.linspace(-3, 3, 128))
25     Z = X**2 - 2*X + X**2 * Y**2 - 2*X*Y
26
27     fig, axs = plt.subplots()
28     co = axs.contourf(X, Y, Z, levels=np.linspace(-10, 100, 80))
29     fig.colorbar(co, ax=axs)
30
31     def overlay_trajectory_contour(axs, trajectory, label, color='k', lw=2):
32         xs=trajectory[:,0]
33         ys=trajectory[:,1]
34         axs.plot(xs,ys, color, label=label, lw=lw)
35         return axs;
36
37     overlay_trajectory_contour(axs, positions, '$\eta = %s' % learning_rate, 'r--', lw=0.5)

```

Mentre per il metodo di Newton si è usato lo stesso codice, sostituendo al learning rate l'inversa della matrice hessiana valutata con gli attuali parametri. Si riporta sotto solo lo spezzone di codice modificato:

```

1     #[...] previous code goes here
2     def update(pos, learning_rate):
3         pos -= learning_rate @ gradient(pos)
4
5     def hess_inv(pos):
6         x = pos[0]
7         y = pos[1]
8         hessian = np.array([[2+2*y**2, 4*x*y-2], [4*x*y-2, 2*x**2]])
9         return np.linalg.inv(hessian)
10
11    #[...] previous code goes here
12    for i in range(iteration_number):
13        positions[i+1] = positions[i]
14        update(positions[i+1], hess_inv(positions[i]))
15
16    #[...] previous code goes here

```

## Esercizio 2 - Metodi di discesa del gradiente

(a) L'esercizio consiste nell'eseguire il fit di dati generati dalla funzione  $\sin(2\pi x)$ , con l'aggiunta di rumore, implementando una regressione polinomiale tramite discesa del gradiente stocastica.

La Figura 2 mostra il risultato del fit usando un polinomio di grado 9 (quindi 10 parametri), un learning rate di 0.2 e un batch size di 20. Questo è stato comparato al risultato ottenuto tramite l'espansione di Taylor (sempre di grado 9) di  $\sin(2\pi x)$  intorno al punto 0.5.

Il fit, nonostante sia graficamente simile al risultato ottenuto con l'espansione, risulta avere dei coefficienti polinomiali che differiscono di ordini di grandezza, e talvolta di segno, da quelli della serie di Taylor.

Si riporta sotto il codice python utilizzato:

```

1     #number of data
2     n_data = 200
3     iteration_number = 100000
4     learning_rate = 0.2

```

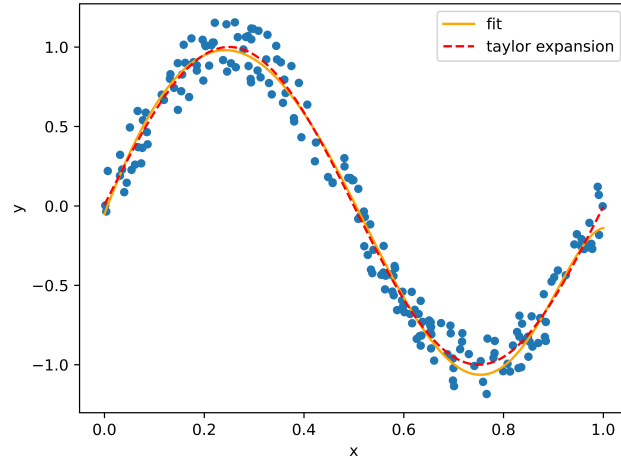


Figura 2: In blu i data points generati, in arancione la curva ottenuta da un fit lineare dei coefficienti di un polinomio di ordine 9, in rosso l'espansione di Taylor intorno a 0.5 (sempre al nono ordine).

```

5 batch_size = 20
6 # pol_order is the order of the polynomial -1
7 pol_order = 10
8
9 x = np.random.rand(n_data)
10
11 y = np.sin(2*np.pi*x)+(np.random.rand(len(x))-0.5)*(0.2/0.5)
12
13 fig, ax = plt.subplots()
14 ax.scatter(x,y, s = 20)
15
16 parameters = np.random.rand(pol_order)-0.5
17
18 def batch(x,y, batch_size):
19     in_batch = np.random.choice(np.arange(len(x)), batch_size)
20     x_batch = np.array([x[i] for i in in_batch])
21     y_batch = np.array([y[i] for i in in_batch])
22     return (np.array(x_batch), np.array(y_batch))
23
24 def polynomial(parameters, x, order):
25     temp = np.zeros(len(x))
26     for i in range(order):
27         temp += parameters[i]*(x**i)
28     return temp
29
30 def gradient(parameters, x, y):
31     temp = np.zeros(len(parameters))
32     pol = polynomial(parameters, x, len(parameters))
33     for i in range(0, len(parameters)):
34         temp[i] = -2*np.sum((y-pol)*(x**i))/len(x)
35     return temp
36
37 def update(parameters, learning_rate, x, y):
38     parameters -= learning_rate * gradient(parameters, x, y)
39
40
41 for i in range(iteration_number):
42     batch_x, batch_y = batch(x,y, batch_size)
43     update(parameters, learning_rate, batch_x, batch_y)
44
45 x_lin = np.linspace(0., 1., n_data)
46 y_predicted = polynomial(parameters, x_lin, pol_order)
47 ax.plot(x_lin, y_predicted, color = 'orange', label = "fit")
48

```

```

49 #comparing with taylor expansion
50 par_taylor = np.array([0.00692527, 6.13253, 1.4848, -50.0794, 34.0311, -10.1117,
51 173.146, \
52 -301.822, 189.264, -42.0587, -3.23628*10**(-15)])
53 y_taylor = polynomial(par_taylor[:pol_order], x_lin, pol_order)
54 ax.plot(x_lin, y_taylor, color = 'red', linestyle = 'dashed', label = "taylor
55 expansion")
56 ax.legend()

```

(b) Si è ottimizzata la funzione:

$$f(x) = x^4 - 3x^3 + 2$$

usando la discesa del gradiente, la discesa del gradiente con momento e la discesa del gradiente accelerata di Nesterov. Un semplice studio della funzione mostra che questa ha il minimo globale in  $x = \frac{9}{4}$  e un punto di sella in  $x = 0$ .

La Tabella 1 mostra la performance dei vari algoritmi al variare dei parametri. Nei casi riportati nella prima e nelle ultime due righe della tabella, la prestazione dei due algoritmi più sofisticati risulta migliore di quella ottenuta dalla discesa del gradiente classica. L'algoritmo classico risulta essere più lento poiché considera solamente il gradiente nella posizione attuale (che vicino ai punti critici è prossimo a zero), mentre quelli più avanzati tengono conto della elevata pendenza iniziale.

$x_0$	$\alpha$	$\gamma$	$n$ gradient descent	$n$ momentum	$n$ Nesterov
8	0.001	0.05	212	201	201
8	0.001	0.5	212	82	98
8	0.001	0.8	212	1204 (sella)	2140 (sella)
-8	0.001	0.5	10000+ (sella)	5318 (sella)	5400 (sella)
-8	0.001	0.8	10000+ (sella)	129	55

Tabella 1: La tabella mostra la performance dei vari algoritmi al variare dei parametri rilevanti.  $x_0$  indica la posizione iniziale,  $\alpha$  il tasso di apprendimento,  $\gamma$  il momentum parameter, ed  $n$  il numero di iterazioni che impiega un determinato algoritmo a raggiungere la destinazione. In caso quest'ultima non fosse il minimo ma il punto di sella, al valore numerico di  $n$  viene accostato un (sella). Se al valore numerico è invece accostato il simbolo +, allora l'algoritmo raggiunge la destinazione, ma si avvicina talmente lentamente da non raggiungerla nel numero di iterazioni che compie.

Nell'ultimo caso riportato in Tabella 1 l'algoritmo classico raggiunge il punto di sella, mentre i due algoritmi più avanzati superano questo per poi oscillare intorno al minimo globale. Questo comportamento è mostrato in Figura 3, ed è sempre spiegabile considerando la "memoria" che gli algoritmi più avanzati possiedono della derivata nella posizione iniziale.

L'algoritmo classico invece mostra una migliore performance nel caso riportato nella terza riga della Tabella 1. Stavolta gli algoritmi più avanzati saltano il minimo globale, per invece oscillare intorno al punto di sella.

L'implementazione dei vari metodi di discesa del gradiente è riportata sotto:

```

1 learning_rate = 0.001
2 momentum_parameter = 0.8
3 iteration_number = 10000
4 initial_pos = -8
5 initial_speed = 0
6
7 #gradient descent method
8 pos_grad = np.zeros(iteration_number+1)
9
10 #gradient descent with momentum method
11 pos_mom = np.zeros(iteration_number+1)
12 speed_mom = np.zeros(iteration_number+1)
13
14 #Nesterov method
15 pos_nesterov = np.zeros(iteration_number+1)

```

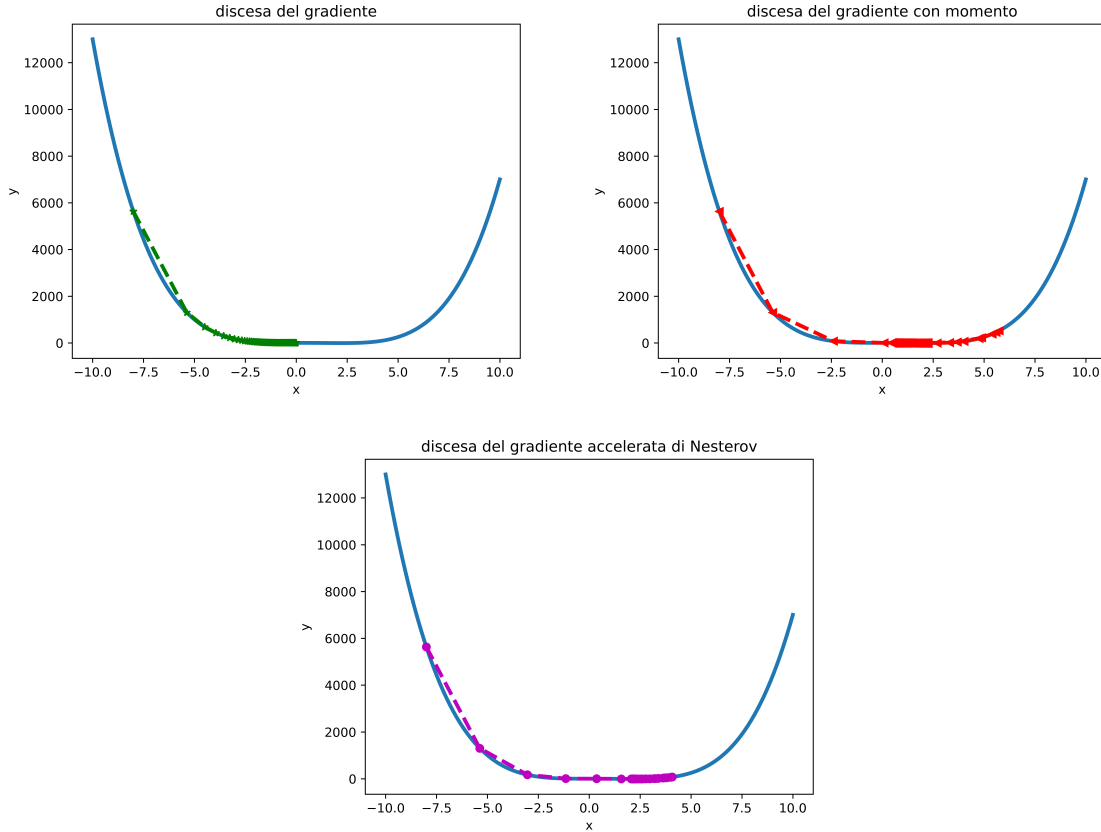


Figura 3: Tre grafici della funzione  $f(x)$ , tracciate su queste le traiettorie corrispondenti ai diversi algoritmi implementati, con tasso di apprendimento  $\alpha = 0.001$ , momentum parameter  $\gamma = 0.8$  e posizione iniziale  $x_0 = -8$ . Osservare che mentre il primo algoritmo non superano il punto di sella, i due algoritmi più sofisticati lo sorpassano, e oscillano intorno al punto di minimo fino a raggiungerlo.

```

16 speed_nesterov = np.zeros(iteration_number+1)
17
18 #setting the initial position and speeds
19 pos_grad[0] = initial_pos
20
21 pos_mom[0] = initial_pos
22 speed_mom[0] = initial_speed
23
24 pos_nesterov[0] = initial_pos
25 speed_nesterov[0] = initial_speed
26
27 #function definitions
28 def gradient(pos):
29     return 4*(pos**3) - 9*(pos**2)
30
31 def update_grad(pos, learning_rate):
32     pos -= learning_rate * gradient(pos)
33     return pos
34
35 def update_mom(pos, speed, learning_rate, momentum_parameter):
36     speed = speed*momentum_parameter + learning_rate*gradient(pos)
37     pos -= speed
38     return (pos, speed)
39
40 def update_nesterov(pos, speed, learning_rate, momentum_parameter):
41     speed = speed*momentum_parameter + learning_rate*gradient(pos-momentum_parameter*
42     speed)
43     pos -= speed
44     return (pos, speed)

```

```

44
45 #booleans to indicate if it has found the minima:
46 reached_destination_grad = False
47 reached_destination_mom = False
48 reached_destination_nesterov = False
49
50 #update loop
51 for i in range(iteration_number):
52     pos_grad[i+1] = update_grad(pos_grad[i], learning_rate)
53     pos_mom[i+1], speed_mom[i+1] = update_mom(pos_mom[i], speed_mom[i], learning_rate,
54         momentum_parameter)
55     pos_nesterov[i+1], speed_nesterov[i+1] = update_nesterov(pos_nesterov[i],
56         speed_nesterov[i], learning_rate, momentum_parameter)
57
58     if not reached_destination_grad:
59         if np.isclose(pos_grad[i+1], 9./4., atol=0.01) and np.isclose(pos_grad[i],
60             9./4., atol=0.01):
61             print("classical gradient descent has reached global minima in ", i, "
62                 iterations")
63             reached_destination_grad = True
64         if np.isclose(pos_grad[i+1], 0., atol=0.01) and np.isclose(pos_grad[i], 0.,
65             atol=0.01):
66             print("classical gradient descent has reached saddle point in ", i, "
67                 iterations")
68             reached_destination_grad = True
69
70     if not reached_destination_mom:
71         if np.isclose(pos_mom[i+1], 9./4., atol=0.01) and np.isclose(pos_mom[i],
72             9./4., atol=0.01):
73             print("gradient descent with momentum has reached global minima in ", i, "
74                 iterations")
75             reached_destination_mom = True
76         if np.isclose(pos_mom[i+1], 0., atol=0.01) and np.isclose(pos_mom[i], 0., atol
77             =0.01):
78             print("gradient descent with momentum has reached saddle point in ", i, "
79                 iterations")
80             reached_destination_mom = True
81
82     if not reached_destination_nesterov:
83         if np.isclose(pos_nesterov[i+1], 9./4., atol=0.01) and np.isclose(pos_nesterov
84             [i], 9./4., atol=0.01):
85             print("nesterov has reached global minima in ", i, "iterations")
86             reached_destination_nesterov = True
87         if np.isclose(pos_nesterov[i+1], 0., atol=0.01) and np.isclose(pos_nesterov[i
88             ], 0., atol=0.01):
89             print("nesterov has reached saddle point in ", i, "iterations")
90             reached_destination_nesterov = True
91
92 #plotting
93 def overlay_trajectory(ax, xs, label, color='k', lw=2):
94     ys = xs**4-3*xs**3+2
95     ax.plot(xs,ys, color, label=label,lw=lw)
96     return ax
97
98 fig1, ax1 = plt.subplots()
99 fig2, ax2 = plt.subplots()
100 fig3, ax3 = plt.subplots()
101
102 x = np.linspace(-10, 10, 128)
103 y = x**4 -3*x**3 +2
104
105 ax1.plot(x,y, lw= 3)
106 ax1.set_xlabel('x')
107 ax1.set_ylabel('y')
108 ax1.set_title("discesa del gradiente")
109
110 ax2.plot(x,y, lw= 3)
111 ax2.set_xlabel('x')
112 ax2.set_ylabel('y')

```

```

103 ax2.set_title("discesa del gradiente con momento")
104
105 ax3.plot(x,y, lw= 3)
106 ax3.set_xlabel('x')
107 ax3.set_ylabel('y')
108 ax3.set_title("discesa del gradiente accelerata di Nesterov")
109
110 overlay_trajectory(ax1,pos_grad,'gradient','g--*', lw=3)
111 overlay_trajectory(ax2,pos_mom,'momentum','r--<', lw=3)
112 overlay_trajectory(ax3,pos_nesterov,'nesterov','m--o', lw=3)

```