

## 13 Quarta Lezione A

### 13.1 Reti neurali

Essendo il computer molto stupido a qualcuno è venuta la brillante idea di cercarlo di educarlo, tanto per divertirsi un po'. Cominciamo quindi questa *descensio ad inferos* di cui al più faremo i primi gradini. Vogliamo infatti costruire un semplice esempio di rete neurale, che serva a classificare dei dati. Quello che vogliamo fare è costruire una scatola nera che prenda in input un certo set di dati a cui corrisponde un certo label, e la nostra scatola nera deve imparare a capire qual è il label a seconda di varie caratteristiche presenti nei dati in input. Metaforicamente potremmo dire che "allenare" una rete neurale è come fare esercizi guardando le soluzioni finché non diventi bravo e impari a farli da solo. Si tratta quindi di apprendimento supervisionato. Vediamo uno schema di una rete neurale e cerchiamo di capire cosa succede.

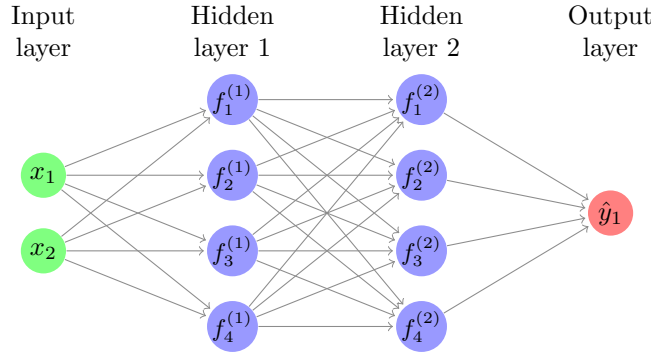


Figura 7: Schema di una rete neurale per una classificazione binaria, eventualmente il layer di output può avere due neuroni in caso di one-hot-encoding

Questo è un piccolo esempio della struttura di una rete neurale. Essa è suddivisa in layer e ogni layer contiene un certo numero di neuroni (i pallini). Il layer di input è quello a cui noi passiamo i dati, mentre il layer di output è quello che ci restituisce il risultato; è negli hidden layer che avviene la magia. Supponiamo di avere dei dati del tipo:

$x_1$	$x_2$	$y$
0.34	0.56	1
0.5	0.89	1
0.2	0.7	0
0.52	0.1	1
0.9	0.83	0
$\vdots$	$\vdots$	$\vdots$

Tabella 1: Tabella di dati da classificare: ogni dato ha due features ovvero le due coordinate  $x_1$  e  $x_2$  ed un label o target ovvero  $y$ , cioè abbiamo assegnato ad ogni punto sul piano un valore binario 0 o 1.

Quindi abbiamo dei punti sul piano, senza perdere di generalità supponiamoli fra 0 ed 1, e li abbiamo targhettati con uno zero o con un uno. Vediamo ora come addestrare la rete. Come prima cosa passiamo i dati in input alla rete e facciamo il primo passo, cioè dobbiamo andare verso il primo layer nascosto, quello che in genere si fa è la seguente cosa:

$$\mathbf{z} = W\mathbf{x} + \mathbf{b} \quad . \quad (81)$$

Non molto magico eh? vediamo gli elementi di questa formula  $\mathbf{x}$  è il vettore in input quindi  $\mathbf{x} = (x_1, x_2)$ ;  $W$  è una matrice e gli elementi vengono chiamati pesi; infine  $\mathbf{b}$  è un vettore ed è chiamato bias. Le dimensioni di questa matrice e questo vettore sono dati dalla dimensione del layer: nella fattispecie per il primo layer  $W$  sarà una matrice  $4 \times 2$  e  $\mathbf{b}$  un vettore lungo 4. Per cui anche il nostro output  $\mathbf{z}$  sarà un vettore lungo 4. Più in generale possiamo dire, essendo questa struttura la stessa per ogni layer, che la dimensione di  $W$  sarà sempre della forma: (numero di neuroni nel layer attuale)  $\times$  (numero di neuroni layer precedente), mentre  $\mathbf{b}$  sarà sempre un vettore lungo (numero di neuroni nel layer attuale). La domanda sorge spontanea: che valori hanno le entrate di  $W$  e di  $\mathbf{b}$ ? Inizialmente saranno scelte in maniera casuale e poi vedremo che durante le iterazioni di allenamento, dette epoche, esse verranno aggiornate secondo una certa regola. Prima di passare al secondo layer però c'è un'altra operazione da fare. Infatti solitamente non è  $\mathbf{z}$  l'output ma  $f(\mathbf{z})$  dove  $f$  è una certa funzione  $f: \mathbb{R} \rightarrow \mathbb{R}$ . Questa  $f$  è chiamata funzione di attivazione del neurone; essa restituirà nel nostro caso un vettore lungo 4 che sarà l'input per il layer successivo. Possiamo vedere tutto ciò come delle regole per far

accendere o spegnere un neurone. Fatto questo il procedimento si ripete in maniera uguale per ogni layer fino ad arrivare al layer di output. Ora che il passaggio attraverso la rete è stato fatto dobbiamo capire se l'output che genera la rete è simile a quello che noi vogliamo. È dunque arrivato il momento di calcolare una funzione che misuri la distanza tra risultato esatto e risultato della rete. Questa funzione è chiamata *loss* e nel nostro caso di classificazione binaria useremo quella che è chiamata binary cross entropy:

$$\mathcal{L} = \frac{1}{N} \sum_i^N (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)) \quad . \quad (82)$$

Non vi sarà sfuggito il fatto che questa è la classica espressione per l'entropia di un sistema di fermioni. Precisiamo che  $N$  è il numero di dati che passiamo per farlo allenare (volendo il numero di righe della tabella di sopra). Fatto questo non ci resta che decidere il modo il cui aggiornare i pesi e bias della rete. Abbiamo detto che la loss misura la distanza di quanto la rete sbaglia, quindi basterà minimizzarla. Trovando il punto di minimo avremo i parametri (pesi e bias) ottimali della nostra rete. Come se stessimo eseguendo un fit, circa. Per farlo usiamo il più semplice e classico algoritmo: il gradiente discendente (di cui è presente una discussione in una delle appendici). Abbiamo quindi che:

$$W^{i+1} = W^i - \alpha \frac{\partial \mathcal{L}}{\partial W^i}, \quad (83)$$

$$\mathbf{b}^{i+1} = \mathbf{b}^i - \alpha \frac{\partial \mathcal{L}}{\partial \mathbf{b}^i}, \quad (84)$$

dove  $\alpha$  è chiamato "learning rate" ed è un parametro che scegliamo noi. Immagino non vi sorprenderà sapere che le derivate sopra citate si calcolano con la chain rule:

$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial f}{\partial z} \frac{\partial (W\mathbf{x} + \mathbf{b})}{\partial W}, \quad (85)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial f}{\partial z} \frac{\partial (W\mathbf{x} + \mathbf{b})}{\partial \mathbf{b}}. \quad (86)$$

Fatto questo abbiamo dei nuovi pesi e bias con cui ripartire prendendo nuovamente i dati in input e propagarli lungo la rete. Questo passaggio si chiama "feed forward" mentre l'updating è chiamato "backpropagation". Fatti questi, è finita un'iterazione, ovvero un'epoca, e parte la successiva. Il numero di epoche è anch'esso un parametro che scegliamo noi. Altra cosa che sta al nostro giudizio è la scelta delle funzioni di attivazione, per le quali, eccezion fatta per il layer di output, non ci sono chissà che criteri per selezionarle. In letteratura si possono trovare le scelte più disparate. Noi useremo per tutti i layer nascosti delle tangenti iperboliche e per il layer di output una sigmoide. Aggiungiamo altre due cose: avendo usato come loss la binary cross entropy può essere interessante misurare l'accuratezza del risultato della rete anche in un altro modo, calcoliamo:

$$a = 1 - \left| \sum_i^N \frac{\hat{y}_i - y_i}{N} \right|, \quad (87)$$

detta accuratezza più siamo vicini ad uno meglio è però bisogna stare attenti a non overfittare. La rete potrebbe infatti star "imparando a memoria" e questo non va bene. Per controllare che questo non succeda usiamo quella è chiamata validation loss. Si tratta sempre di calcolare la loss, ma non sui dati che usiamo come dati di allenamento, ma su un altro set che usiamo esclusivamente per questo calcolo. Quindi noi passiamo alla rete dei dati e li dividiamo in due, su un set allena e sull'altro valida. Quello che facciamo è plottare insieme la loss calcolata sui dati di allenamento e quella calcolata sui dati di validation (in funzione delle epoche). Se le due scendono insieme la rete si sta comportando bene. Se invece ad un certo punto la loss scende ma la loss di validation sale vuol dire che la rete sta overfittando, cioè sta imparando a memoria le caratteristiche dei dati di allenamento. Quindi per evitare questo bisogna scegliere alcuni parametri con un po' di senno, ad esempio learning rate, numero di neuroni e numero di layer o anche numero di epoche. Detto questo passiamo ora a vedere il codice. Implementeremo una rete neurale con un solo layer nascosto con numero di neuroni variabile, due neuroni in input e uno solo in output. Lo scopo sarà proprio di riconoscere dati dei punti nel quadrato  $[0, 1] \times [0, 1]$  se essi hanno un valore associato 1 oppure 0. L'inizializzazione verrà fatta estraendo dei numeri distribuiti gaussianamente. Omettiamo per brevità, come sempre, le funzioni che fanno i plot e quello che è il main del codice.

```
1 """
2 Code that implement a shallow neural network for a binary classifications.
3 The code is witten imposed 2 input, 1 output e one hidden layer.
4 Is possible to choose the dimesions of hidden layer.
5 It is also possible to save plots during the run to see how the network is learning.
```

```

6 """
7 import numpy as np
8 import matplotlib.pyplot as plt
9
10 #=====
11 # Loss function binary classification
12 #=====
13
14 def Loss(Yp, Y):
15     '''
16     loss function, binary crosss entropy
17
18     Parameters
19     -----
20     Yp : 1darray
21         actual prediction
22     Y : 1darray
23         Target
24
25     Returns
26     -----
27     float, binary crosss entropy
28     '''
29     m = len(Y)
30     return -np.sum(Y*np.log(Yp) + (1 - Y)*np.log(1 - Yp))/m
31
32 #=====
33 # Activation function
34 #=====
35
36 # Hidden layer
37 def g1(x):
38     return np.tanh(x)
39 # Output layer
40 def g2(x):
41     return 1 / (1 + np.exp(-x))
42
43 #=====
44 # Initialization
45 #=====
46
47 def init(n):
48     '''
49     Random initialization of parameters weights and biases
50
51     Parameters
52     -----
53     n : int
54         number of neurons in the hidden layer
55
56     Returns
57     -----
58     W1, b1 : 2darray
59         weights and bias for hidden layer
60     W2, b2 : 2darray
61         weights and bias for output layer
62     '''
63     # Hidden layer
64     # nx2 because 2 featurss and n neurons
65     W1 = np.random.randn(n, 2)
66     b1 = np.random.rand(n, 1)
67     # Output layer
68     # 1xn because 1 output and n neurons
69     W2 = np.random.randn(1, n)
70     b2 = np.random.rand(1, 1)
71     return W1, b1, W2, b2
72
73 #=====
74 # Network prediction function
75 #=====
76
77 def predict(X, W1, b1, W2, b2):
78     """
79     Function that returns the prediction of the network
80
81     Parameters

```

```

82  -----
83  X : 2darray
84      data, features
85  W1, b1, W2, b2 : 2darray
86      parameter of the network
87
88  Returns
89  -----
90  A1 : 1d array
91      intermediate prediction
92  A2 : 1d array
93      final prediction
94  """
95  # Hidden layer
96  Z1 = W1 @ X + b1
97  A1 = g1(Z1)
98  # Output layer
99  Z2 = W2 @ A1 + b2
100  A2 = g2(Z2)
101
102  return A1, A2
103
104  #=====
105  # Backpropagation function
106  #=====
107
108  def backpropagation(X, Y, step, A1, A2, W1, b1, W2, b2):
109      """
110      Backpropagation function.
111      Update weights and biases with gradient descent
112      all the quantities came from taking the derivative of the Loss
113
114      Y : 1darray
115          Target
116      step : float
117          learning rate
118      A1, A2 : 1darray
119          predictions of the network
120      W1, b1, W2, b2 : 2darray
121          parameter of the network
122      """
123      m = len(Y)
124      # Output layer
125      dLdZ2 = (A2 - Y)
126      dLdW2 = dLdZ2 @ A1.T / m
127      dLdb2 = np.sum(dLdZ2, axis=1)[:, None] / m
128      # Hidden layer
129      dLdZ1 = W2.T @ dLdZ2 * (1 - A1**2)
130      dLdW1 = dLdZ1 @ X.T / m
131      dLdb1 = np.sum(dLdZ1, axis=1)[:, None] / m
132
133      # Update of parameters
134      W1 -= step * dLdW1
135      b1 -= step * dLdb1
136      W2 -= step * dLdW2
137      b2 -= step * dLdb2
138
139      return W1, b1, W2, b2
140
141  #=====
142  # Accuracy mesuraments
143  #=====
144
145  def accuracy(Yp, Y):
146      """
147      accuracy of prediction. We use:
148      accuracy = 1 - | sum ( prediction - target )/target_size |
149
150      Parameters
151      -----
152      Yp : 1darray
153          actual prediction
154      Y : 1darray
155          Target
156
157      Returns

```

```

158     -----
159     a : float
160         accuracy
161     '''
162     m = len(Y)
163     a = 1 - abs(np.sum(Yp.ravel() - Y)/m)
164     return a
165
166     =====
167     # Train of the network
168     =====
169
170 def train(X, Y, n_epoch, neuro, step, sp=False, verbose=True):
171     '''
172     function for the training of the network
173
174     Parameters
175     -----
176     X : 2darray
177         data, featur
178     Y : 1darray
179         Target
180     n_epoch : int
181         number of epoch
182     neuro : int
183         number of neurons in the hidden layer
184     step : float
185         learning rate
186     sp : boolean, optional, default False
187         if True a plot of boundary is saved each 100 epoch
188         usefull for animations
189     verbose : boolean, optional, default True
190         if True print loss and accuracy each 100 epoch
191
192     Returns
193     -----
194     result : dict
195         params -> W1, b1, W2, b2 weights and bias of network
196         train_Loss -> loss on train data
197         valid_Loss -> loss on validation data
198     '''
199
200     W1, b1, W2, b2 = init(neuro)
201     L_t = np.zeros(n_epoch) # training loss
202     L_v = np.zeros(n_epoch) # validation loss
203     N = X.shape[1]          # total number of data
204     M = N//4                # nuber of data for validation
205
206     # split dataset in validation and train
207     X_train, Y_train = X[:, :N-M ], Y[:N-M ]
208     X_valid, Y_valid = X[:, N-M:], Y[ N-M:]
209
210     for i in range(n_epoch):
211         # train
212         A1, A2 = predict(X_train, W1, b1, W2, b2)
213         L_t[i] = Loss(A2, Y_train)
214         # validation
215         _, Yp = predict(X_valid, W1, b1, W2, b2)
216         L_v[i] = Loss(Yp, Y_valid)
217         # update
218         W1, b1, W2, b2 = backpropagation(X_train, Y_train, step, A1, A2, W1, b1, W2, b2)
219
220         if not i % 100:
221             #if sp : plot(X_train, Y_train, (W1, b1, W2, b2), i)
222
223             if verbose:
224                 acc = accuracy(A2, Y_train)
225                 print(f'Loss = {L_t[i]:.5f}, accuracy = {acc:.5f}, epoch = {i} \r', end='')
226
227     if verbose: print()
228
229     result = {'params'      : (W1, b1, W2, b2),
230              'train_Loss'  : L_t,
231              'valid_Loss'  : L_v,
232              }
233

```

Dopo queste belle tre paginate di codice vediamo un po' di risultati. Questo codice era scritto per far vedere come rete impara creando un gif, che trovate disponibile sulla cartella. Il risultato finale è il seguente:

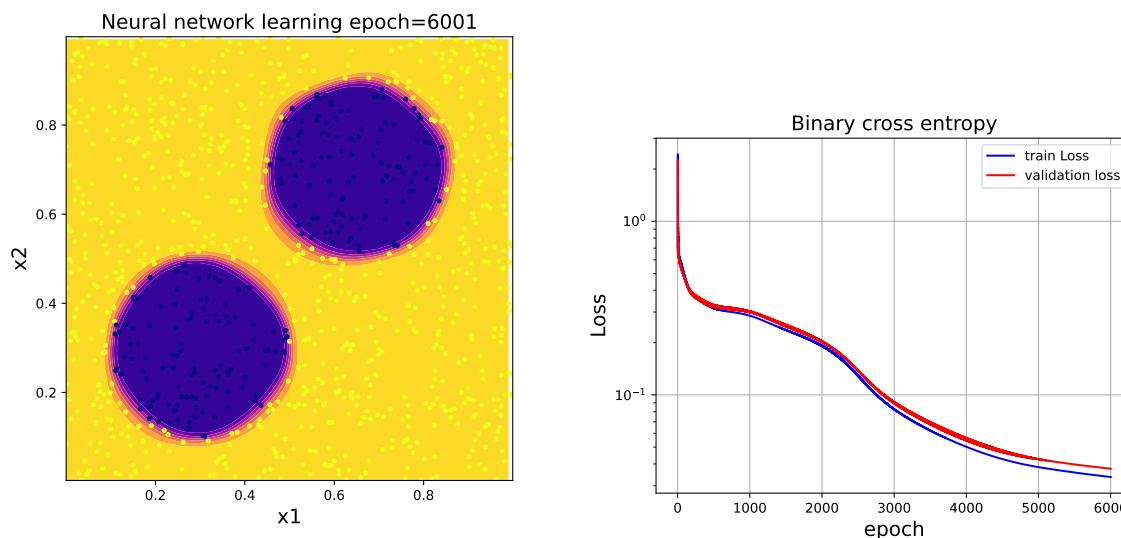


Figura 8: Risultati della rete neurale, predizione e andamento della loss. I parametri qui sono: 20 neuroni nel layer nascosto, un learning rate di 1.5, 6001 epoche, 3000 dati di train e 1000 di validation.

Il plot a sinistra è fatto sui dati di test, quindi un set di dati che la rete non ha usato per allenarsi e il risultato sembra soddisfacente, vediamo poi che le due loss scendono insieme il che è segno che la rete si comporta bene. Volendo ora farvi vedere andamenti diversi della loss vi mostro due grafici provenienti da un altro codice, in cui ho implementato una rete neurale in maniera leggermente diversa e un po' più generale. Anche perché, se vedete la loss nella figura sopra, noterete che la linea è un po' spessa, dovuto al fatto che sta oscillando (immagino sia una concausa tra problema da affrontare e algoritmo di minimizzazione). Il problema è questa volta il MNIST, ovvero il riconoscimento di cifre scritte a mano in bassa risoluzione, ogni immagine è  $28 \times 28$  pixel. Vediamo un caso in cui la rete overfitta e un caso in cui si comporta bene:

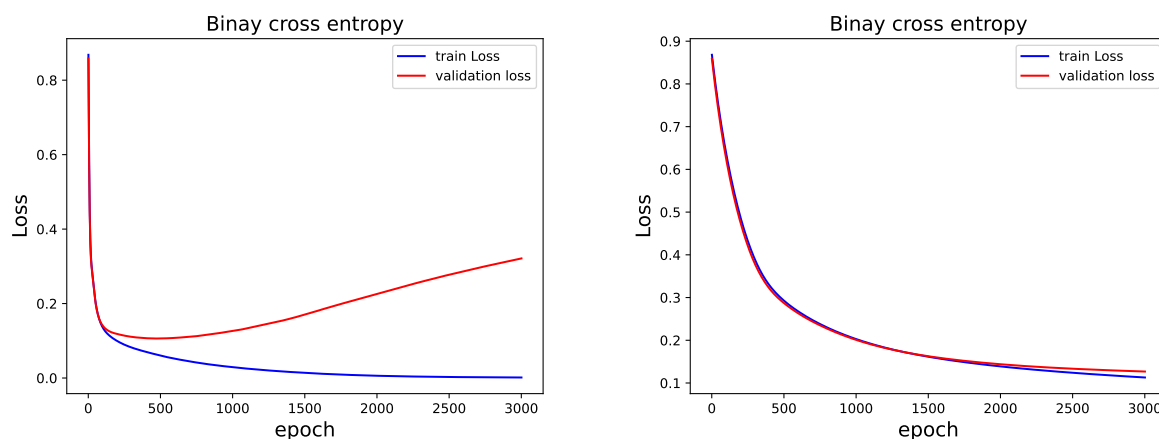


Figura 9: Andamento della loss in funzione delle epoche. Nel grafico a sinistra la rete sta overfittando, mentre a destra l'allenamento procede bene. I parametri qui sono: 2 layer nascosti ognuno da 50 neuroni; l'ottimizzatore usato è ADAM e il learning rate iniziale è di 0.05 per il caso di overfit e di 0.01 per quello di buon fit; le funzioni di attivazione sono: per i layer nascosti le "relu", mentre una sigmoide per il layer finale. Le epoche sono 3001, 2250 i dati di train e 750 quelli di validation.

Vediamo quindi che anche solo il cambiamento del learning rate iniziale (dove si specifica iniziale in quanto ADAM è un algoritmo adattivo) provoca un comportamento che non vogliamo, facendo overfittare la rete. Questo lo possiamo vedere anche, trattandosi di un classificatore, grazie a quelle che sono le matrici di confusione. Mostriamo di seguito le matrici nel caso di overfit e di fit facendo anche il confronto vedendo i dati train

e quelli di test. Spieghiamo velocemente come si legge questa matrice: sull'asse delle ordinate è presente la predizione della rete mentre su quello delle ascisse ci sta il risultato esatto. Le entrate di questa matrice ci dicono fondamentalmente quante volte la rete ha detto "x" e doveva dire "y"; quindi gli elementi sulla diagonale sono le risposte esatte mentre i termini fuori sono risposte sbagliate. Per chiarezza precisiamo che questa matrice viene calcolata sempre a fine allenamento.

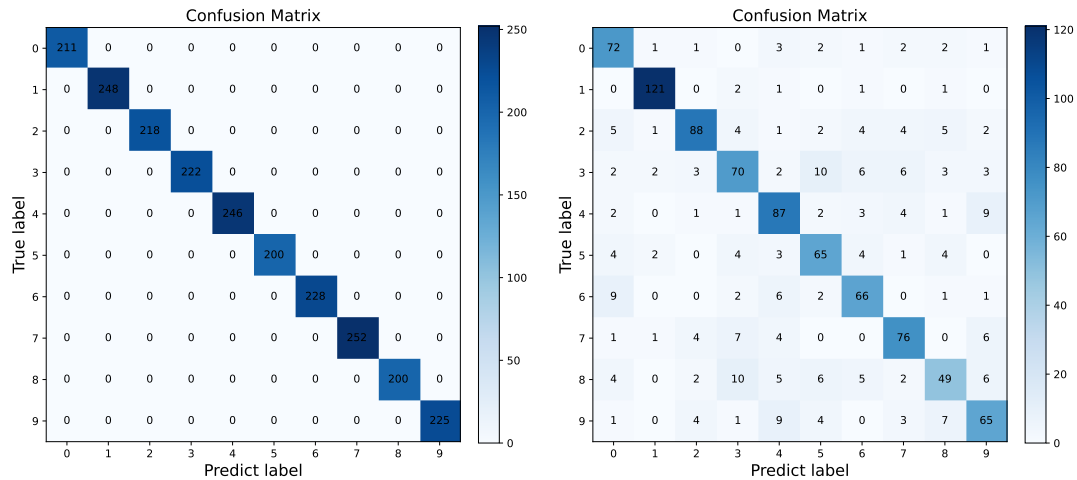


Figura 10: Vediamo qui le matrici nel caso di overfit calcolate sui dati di train, a sinistra, e su quelli di test, a destra. Avendo imparato a memoria le caratteristiche dei dati di train vediamo effettivamente che la rete non sbaglia una singola predizione. Mentre a destra, su dati che per la rete sono nuovi, ci sono molti più errori.

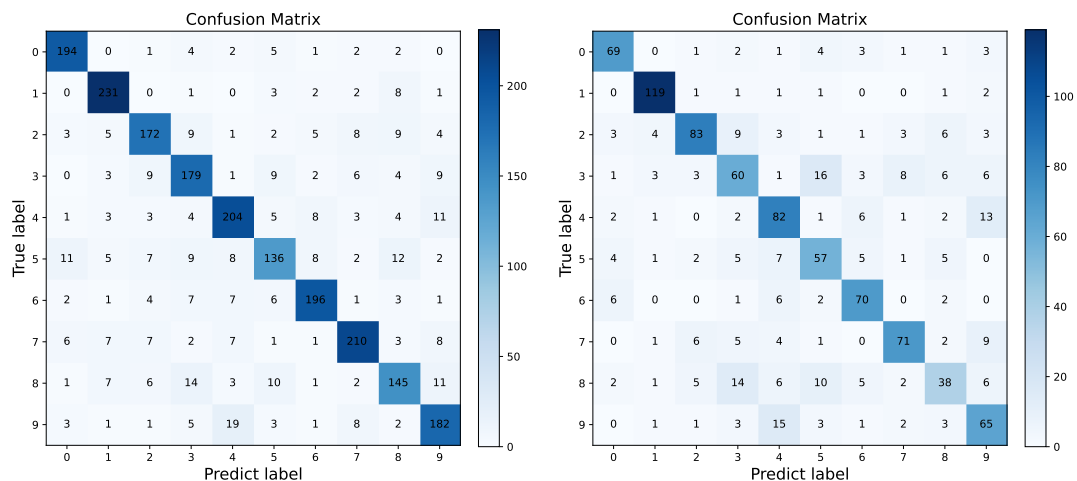


Figura 11: Vediamo qui le matrici nel caso di buon fit calcolate sui dati di train, a sinistra, e su quelli di test, a destra. Questa volta la rete non ha imparato a memoria, vediamo infatti che ci sono risposte sbagliate un po' ovunque. Le due matrici sembrano parenti.

Per concludere questa discussione in caso fosse di vostro interesse vi mostro il codice per calcolare queste matrici di confusione:

```
1 def confmat(true_target, pred_target, plot=True, k=0):
2     """
3     Function for creation and plot of confusion matrix
4
5     Parameters
6     -----
7     true_target : 1darray
8         vaules that must be predict
9     pred_target : 1darray
10        values that the network has predict
```

```

11 plot : bool, optional, default True
12     if True the matrix is plotted.
13 k : int, optional, default 0
14     number of figure, necessary in order not to overlap figures
15
16 Return
17 -----
18 mat : 2darray
19     confusion matrix
20 '''
21
22 dat = np.unique(true_target)      # classes
23 N   = len(dat)                   # Number of classes
24 mat = np.zeros((N, N), dtype=int) # confusion matrix
25
26 # creation of confusion matrix
27 for i in range(len(true_target)):
28     mat[true_target[i]][pred_target[i]] += 1
29
30 if plot :
31     fig = plt.figure(0, figsize=(7, 7))
32     ax = fig.add_subplot()
33
34     c = ax.imshow(mat, cmap=plt.cm.Blues) # plot matrix
35     b = fig.colorbar(c, fraction=0.046, pad=0.04)
36     # write on plot the value of predictions
37     for i in range(mat.shape[0]):
38         for j in range(mat.shape[1]):
39             ax.text(x=j, y=i, s=mat[i, j],
40                     va='center', ha='center')
41
42     # Label
43     ax.set_xticks(dat, dat)
44     ax.set_yticks(dat, dat)
45     ax.tick_params(top=False, bottom=True, labeltop=False, labelbottom=True)
46
47     plt.xlabel('Predict label', fontsize=15)
48     plt.ylabel('True label', fontsize=15)
49     plt.title('Confusion Matrix', fontsize=15)
50     plt.tight_layout()
51
52 return mat

```