

OTTIMIZZAZIONE DI RETI NEURALI ARTIFICIALI: COME SFUGGIRE AI PUNTI DI SELLA?

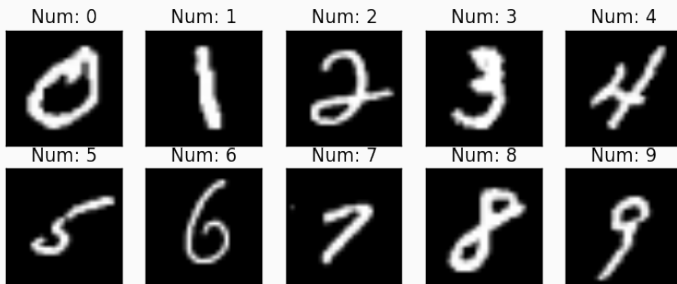
Niccolò Ajroldi
Matricola 846266

Relatore: Danilo Ardagna

25 Febbraio 2019

Politecnico di Milano - Ingegneria Matematica

UN PROBLEMA DI CLASSIFICAZIONE



Come costruire una **funzione** che sappia classificare correttamente le immagini?

Image source: <http://yann.lecun.com/exdb/mnist/>, Yann LeCun, Corinna Cortes, Christopher J.C. Burges

NEURAL NETWORK

Rete di neuroni interconnessi, ciascuno elabora i pixel dell'immagine ed emette un output che ne rappresenta l'attivazione.

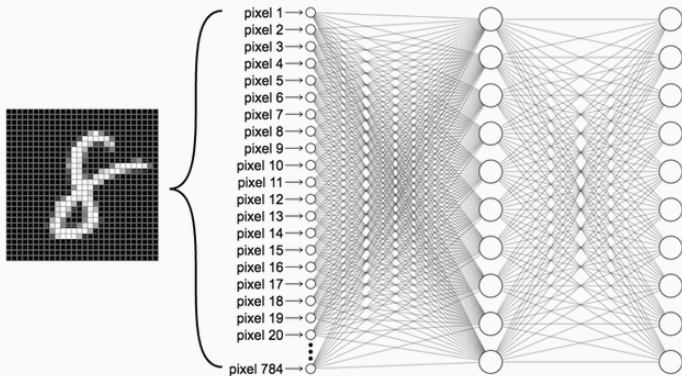


Image source:

<https://achintavarna.wordpress.com/2017/11/17/keras-tutorial-for-beginners-a-simple-neural-network-to-identify-numbers-mnist-data/>

NEURAL NETWORK

Risposta della rete

$$\hat{\underline{y}} = \begin{bmatrix} 0.241 \\ 0.084 \\ 0.316 \\ 0.107 \\ 0.064 \\ 0.381 \\ 0.430 \\ 0.024 \\ 0.879 \\ 0.109 \end{bmatrix}$$

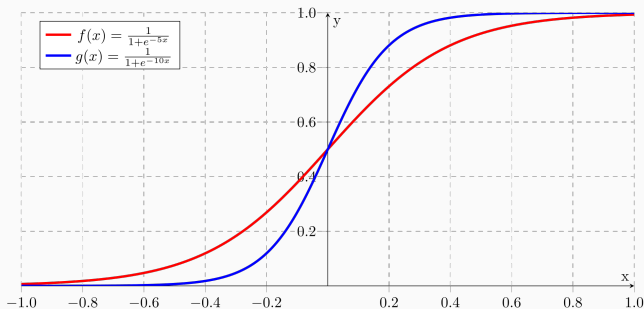
Risposta esatta

$$\underline{y} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

ATTIVAZIONE DI UN NODO

L'output del j -esimo nodo è:

$$z_j = \sigma\left(\sum_{i=1}^{\text{\#input}} w_{ij} * x_i + b_j\right)$$



COME SCEGLIERE I PARAMETRI?

I parametri sono inizializzati random.

Si definisce una funzione che valuti l'**errore** commesso dalla rete neurale nel classificare le N immagini nel dataset:

$$MSE(\underline{w}, \underline{b}) = \frac{1}{N} \sum_{i=1}^N \|\underline{y}_i - \hat{\underline{y}}(x_i, \underline{w}, \underline{b})\|^2$$

Il problema diventa trovare $\underline{w}^*, \underline{b}^*$ tali che:

$$(\underline{w}^*, \underline{b}^*) = \arg \min_{(\underline{w}, \underline{b}) \in \mathbb{R}^{m \times n}} MSE(\underline{w}, \underline{b})$$

⇒ Algoritmi di **Ottimizzazione**

Potenziati problemi:

- Minimi locali
- **Selle**

L'MSE è definito su uno spazio parametrico multidimensionale, nel caso di reti profonde si hanno centinaia di migliaia di parametri.

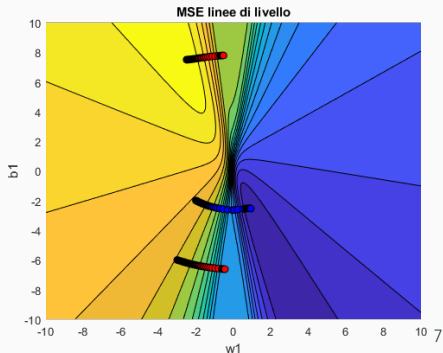
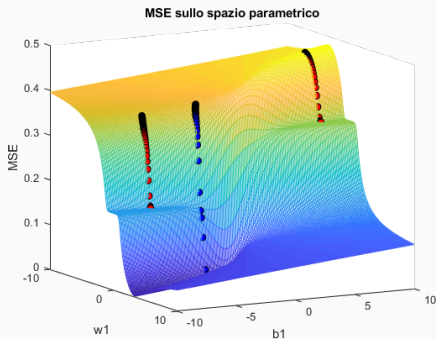
Per visualizzare le dinamiche di ottimizzazione si considera prima un caso molto semplice.

1 NEURONE: 2 PARAMETRI

$$\{(x_i, y_i)\}, x_i \in \mathbb{R}, y_i = \begin{cases} 1 & \text{se } x_i \geq 3 \\ 0 & \text{se } x_i < 3 \end{cases}$$

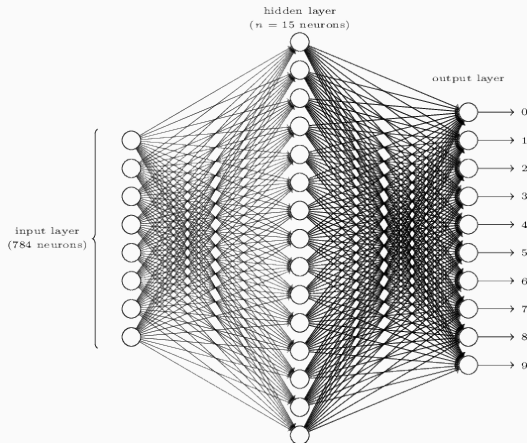
Un unico neurone, la rete è quindi la funzione

$$y = \hat{y}(x, w_1, b_1) = \sigma(w_1 * x + b_1)$$



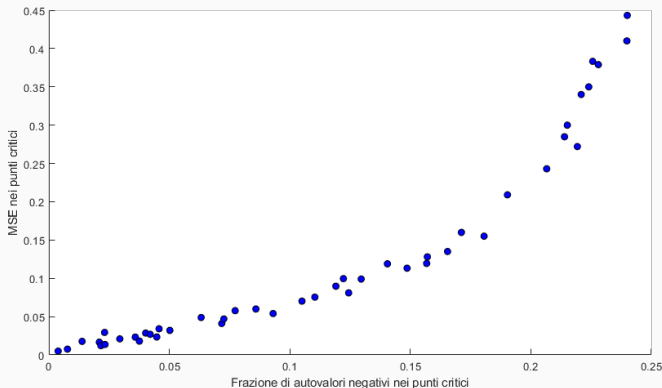
15 NEURONI: 1675 PARAMETRI

Si implementa una rete neurale con **1675 parametri** per la classificazione di immagini.



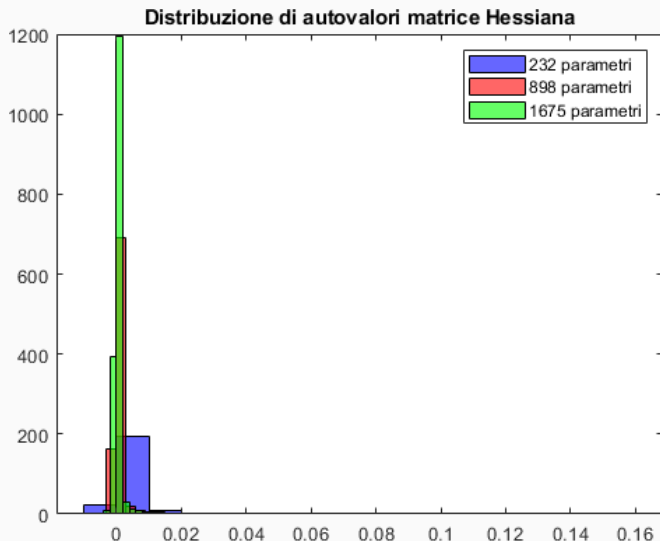
MSE NEI PUNTI ESTREMANTI

Si allena la rete e si valuta l'**Hessiana** nel punto critico raggiunto.



Si osserva una forte correlazione tra la frazione di **autovalori negativi** e l'**errore** nei punti estremanti.

Aumentando le dimensioni della rete lo spettro di autovalori si concentra attorno a zero.



COME FUGGIRE DAI PUNTI DI SELLA?

Metodi del secondo ordine (e.g. Newton)

Si sfrutta l'informazione sulla curvatura contenuta nella matrice hessiana.

- Hessiana malcondizionata
- Elevato costo computazionale

#Neuroni	#Parametri	Tempo
5	565	≈ 1 min
10	1120	≈ 5 min
15	1675	≈ 13 min
20	2230	≈ 25 min

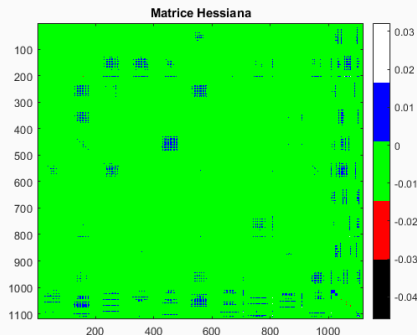
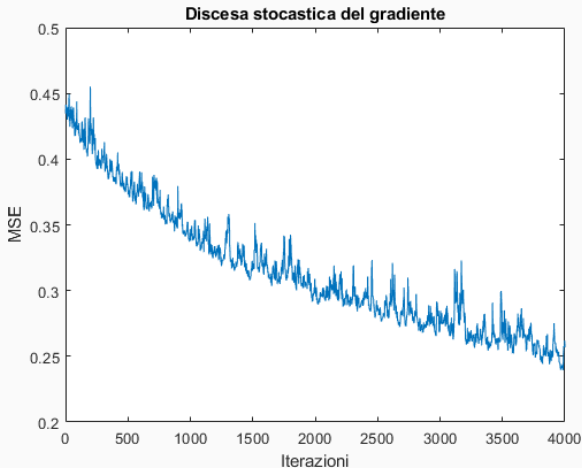


Figure: Hessiana, 1120 parametri

Discesa stocastica del gradiente

Non si ottimizza l'MSE su tutto il dataset ma su una porzione ristretta di esso.

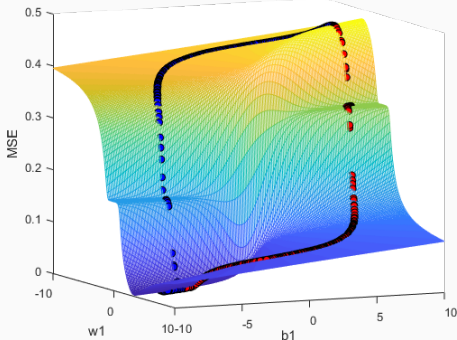
- Semplicità computazionale (primo ordine)
- Può essere implementato in parallelo
- Evita punti critici ad alto errore



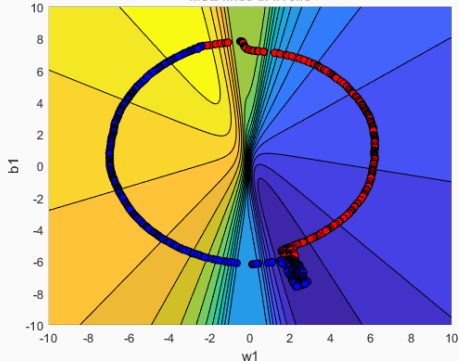
1 NEURONE: 2 PARAMETRI

Si riscontra la capacità del metodo stocastico di superare, o aggirare, i punti di sella in cui il metodo del gradiente si era bloccato.

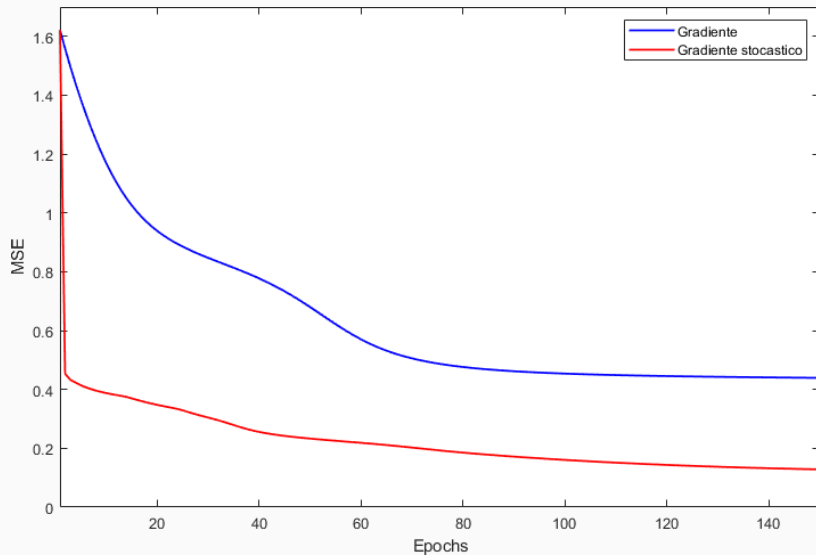
MSE sullo spazio parametrico



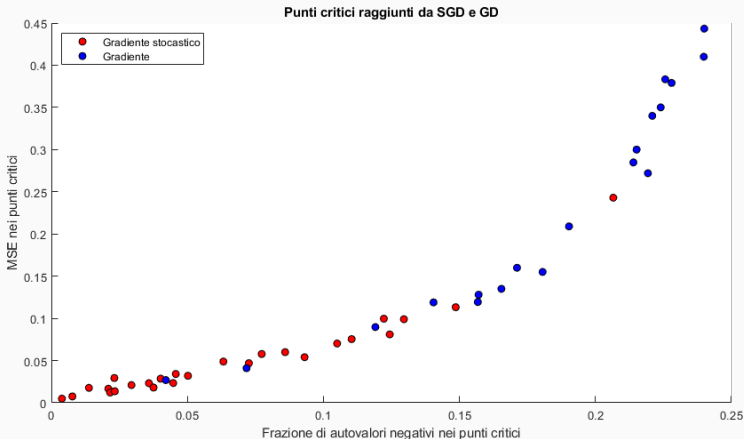
MSE linee di livello



15 NEURONI: 1675 PARAMETRI



Anche nel caso multidimensionale il metodo di discesa stocastica performa decisamente meglio del metodo del gradiente, raggiungendo punti ad errore minore aventi una percentuale inferiore di autovalori negativi.



- Yann N. Dauphin et al. - Identifying and attacking the saddle point problem in high-dimensional non-convex optimization
- Yann LeCun et al. - Eigenvalues of the hessian in deep learning: singularity and beyond
- Catherine F. Higham, Desmond J. Higham - Deep learning: an introduction for applied mathematicians
- Christopher M. Bishop - Exact calculation of the hessian matrix for the multi-layer perceptron