



L'intelligenza biologica (le reti neurali)

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Sommario



Il neurone, modelli deterministici (L-system) e stocastici (frattali).

Reti Neurali.

Mappe topologiche e clustering.

Apprendimento con Rinforzo.

RBF: reti neurali con neuroni a base radiale.

La corteccia

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Le reti neurali



Se il neurone biologico consente l'intelligenza, perché non dovrebbe consentire l'intelligenza artificiale un neurone sintetico?

“.. a neural network is a system composed of *many simple processing elements* operating in *parallel* whose function is determined by *network structure*, *connection strengths*, and the *processing performed at computing elements* or nodes. ...
Neural network architectures are inspired by the architecture of biological nervous systems, which use many simple processing elements operating in parallel to obtain high computation rates”.
(DARPA, 1988)....

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



A cosa servono?



Le reti neurali offrono i seguenti specifici vantaggi nell'elaborazione dell'informazione:

- Apprendimento basato su esempi (non è richiesta l'elaborazione di un modello aderente alla realtà)
- Autoorganizzazione dell'informazione nella rete
- Robustezza ai guasti (codifica ridondante dell'informazione)
- Funzionamento in tempo reale (realizzazione HW)



Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Cosa sono le reti neurali artificiali?



- Le reti neurali sono algoritmi non lineari per l'**approssimazione** di soluzioni di problemi dei quali non esiste un modello preciso (o se esiste è troppo oneroso computazionalmente), mediante l'utilizzo di esempi (dati e uscite) oppure per classificazioni. Connessioni con il dominio della statistica.
- Sono un capitolo importante negli argomenti di intelligenza artificiale.
- Da un altro punto di vista possono essere utilizzate per lo studio delle reti neurali naturali, ovvero dei processi cognitivi

Copyright N.A. Borghese Università di Milano 26/03/2003

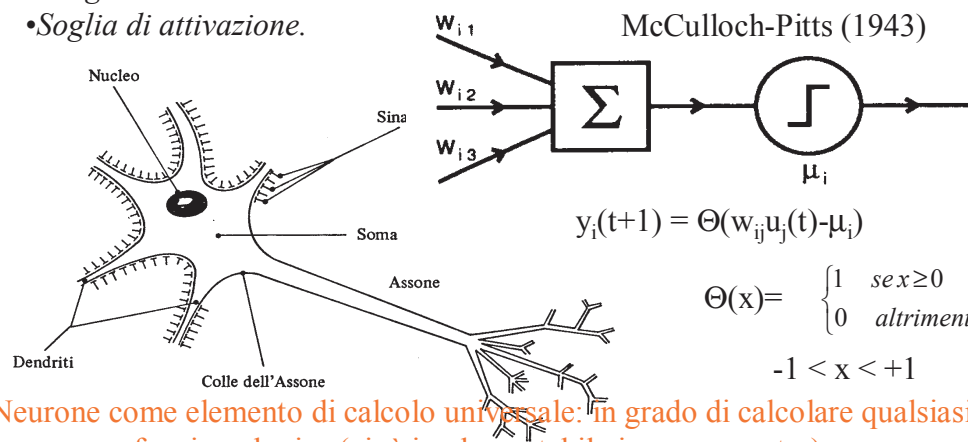
<http://homes.dsi.unimi.it/~borghese>



Il neurone artificiale



- *Potenziale di azione (tutto o nulla).*
- *Integrazione nel soma.*
- *Soglia di attivazione.*



Neurone come elemento di calcolo universale: in grado di calcolare qualsiasi funzione logica (cioè implementabile in un computer).

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Critica al modello di McCulloch-Pitts



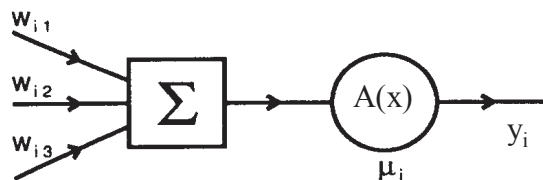
- I neuroni reali non possono essere ridotti ad un dispositivo a soglia. Lo spike ha la sua forma continua che ha una durata di qualche millisecondo.
- Il tempo di propagazione lungo i dendriti non viene considerato.
- La variazione delle forma d'onda del potenziale di membrana lungo il dendrita non viene considerata.
- Gli input non sono sincroni.
- Le interazioni tra input non sono lineari.
- I pesi sono supposti costanti.

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



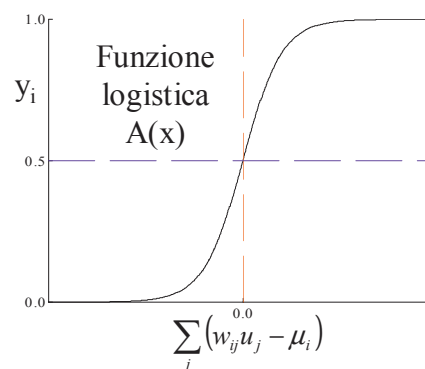
Il perceptrone (Roseblatt, 1962)



Uscita: singolo spike o frequenza di scarica.

Neurone asincrono.
 Soglia μ_i -> traslazione.
 Pesi $\{w_{ij}\}$ -> pendenza.

$$y = g\left(\sum_j (w_{ij}u_j - \mu_i)\right)$$

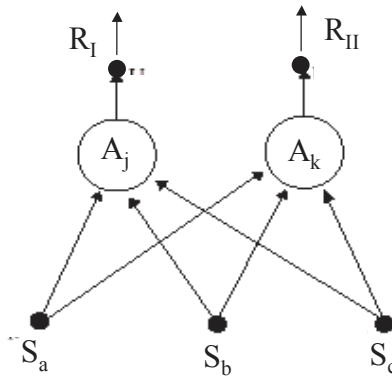


Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Le reti di perceptroni



Apprendimento è la modifica dei parametri in funzione dei parametri di input/output.

$$y = g \left(\sum_j (w_{ij} u_j - \mu_i) \right) = \frac{1}{1 + e^{-\left(\sum_{i=0}^N w_i u_i - \mu \right)}}$$

Questa rete non riesce ad apprendere però funzioni non linearmente separabili quali l'XOR (Minski, 1968).

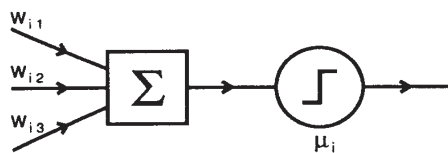


Funzione di attivazione del perceptrone



$$y = \Theta \left(\sum_{j=1} (w_{ij} u_j - \mu_i) \right)$$

$$y = \text{sgn} \left(\sum_{j=1} (w_{ij} u_j - \mu_i) \right)$$



$$y = \text{sgn} \left(\sum_{j=0} (w_{ij} u_j) \right)$$

$$w_{i0} = -\mu_i$$

$$y = \text{sgn}(\mathbf{w} \cdot \mathbf{u})$$

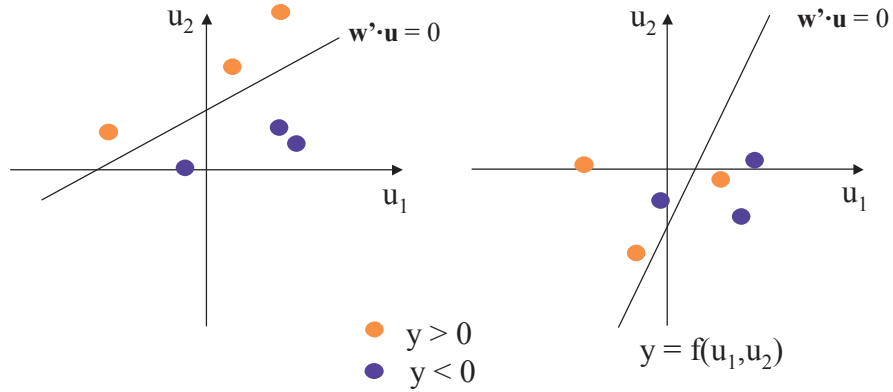


Funzioni linearmente separabili



Linearmente separabile

Non linearmente separabile

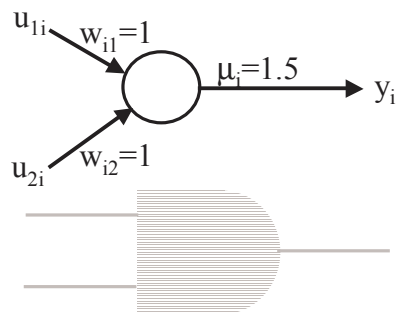


Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Esempio - AND



u_1	u_2	y
-1	-1	-1
-1	1	-1
1	-1	-1
1	1	1

Iperpiano di separazione ($u_0=1$, $w_0 = -\mu_0$):

$$-1.5 + 1 \cdot u_1 + 1 \cdot u_2 = 0$$

$$w_0 u_0 + w_1 u_1 + w_2 u_2 = 0$$

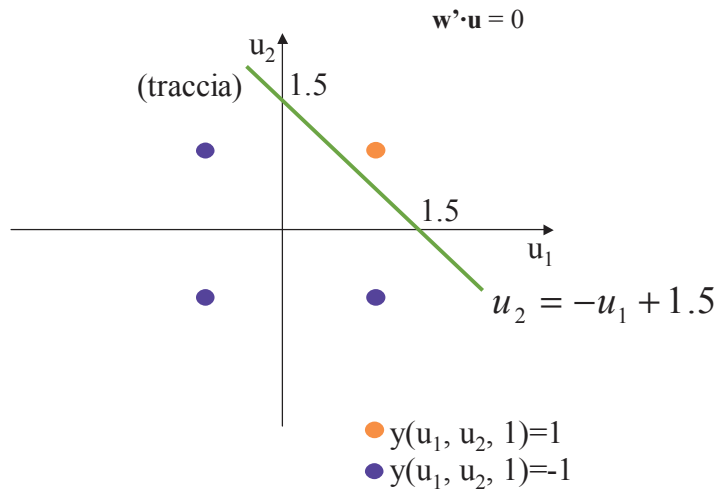
$$u_2 = -u_1 + 1.5$$

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Esempio - AND (grafica)



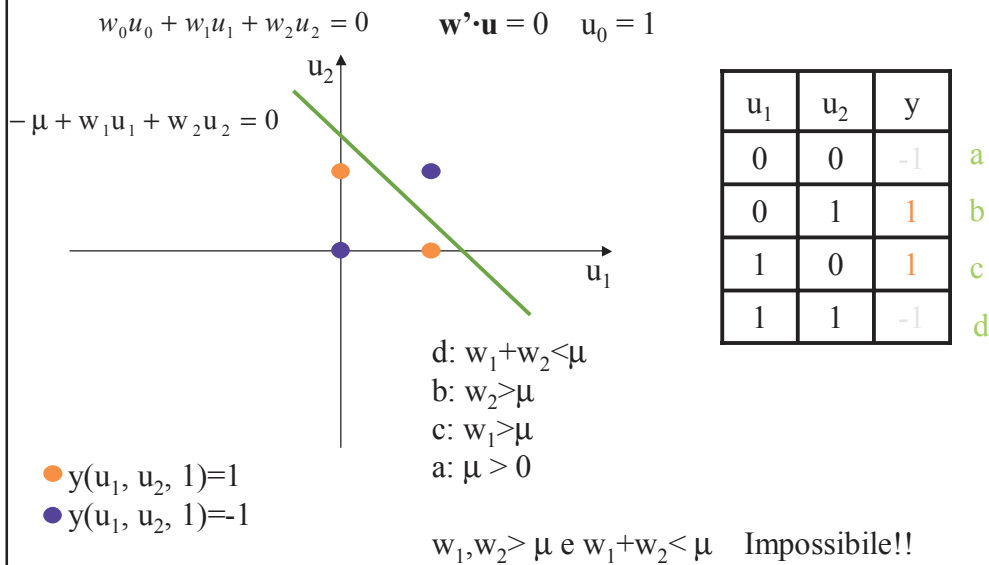
Esistono più soluzioni

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Esempio - XOR

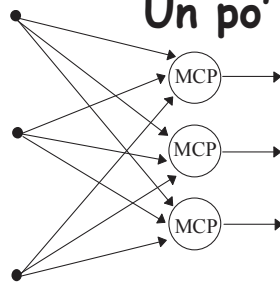


Copyright N.A. Borghese Università di Milano 26/03/2003

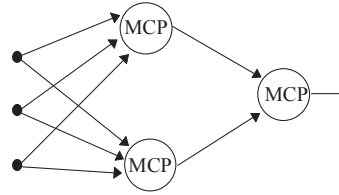
<http://homes.dsi.unimi.it/~borghese>



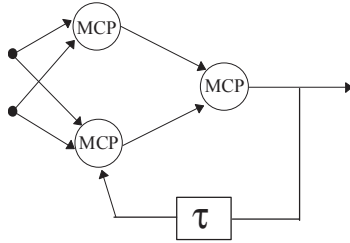
Un po' di tassonomia



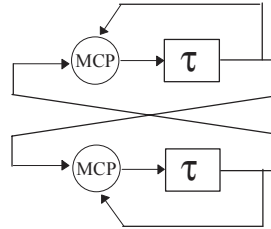
Perceptrone semplice



Perceptrone multistrato



Ricorrente



Ricorrente completamente connessa: autoassociativa (ingresso=stato)

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Complessità della funzione realizzabile

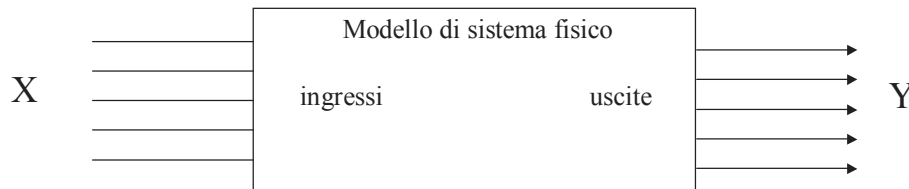


Quanti più neuroni artificiali vengono connessi tanto più la funzione complessiva approssimabile diviene più complessa

$$Y = |f(y_1), f(y_2), f(y_3), \dots, f(y_n)|^T$$

$$y_i = g(X)$$

$$X = |x_1, x_2, x_3, \dots, x_n|^T$$



Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Riassunto



I neuroni connessioneisti sono basati su:

- Ricevere una somma pesata degli ingressi.
- Trasformarla secondo una funzione non-lineare (scalino o logistica)
- Inviare il risultato di questa funzione all'uscita o ad altre unita'.

Le reti neurali sono topologie ottenute connettendo tra loro i neuroni in modo opportuno e riescono a calcolare funzioni molto complesse.



I vari tipi di apprendimento



Supervisionato (learning with a teacher). Viene specificato per ogni pattern di input, il pattern desiderato in input.

Non-supervisionato (learning without a teacher). I neuroni verranno associati a pattern di ingresso contigui. Clustering. Mappe neurali.

Apprendimento con rinforzo (reinforcement learning, learning with a distal teacher). L'ambiente fornisce un'informazione del tipo success or fail.



Apprendimento supervisionato



- Obiettivo: se esiste una soluzione, trovare ΔW in modo iterativo tale che l'insieme dei pesi W^{nuovo} ottenuto come:

$$W^{\text{nuovo}} = W^{\text{vecchio}} + \Delta W$$

dia luogo a un errore sulle uscite minore di W^{vecchio}

- Si tratta di un problema di minimizzazione di una cifra di merito (J) sullo spazio di parametri W .

$$\min_{\{w\}} J(\cdot)$$

$$J = \|Y^D - g(W^{\text{nuovo}}U)\| \leq \|Y^D - g(W^{\text{vecchio}}U)\|$$

Y^D è l'uscita desiderata nota.

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Apprendimento supervisionato



Coppie input/output note.

Definizione di una funzione costo che misuri l'errore sull'uscita.

Modifica dei valori dei pesi in modo tale che la funzione costo sia minimizzata.

Reti multi-strato hanno elevata capacità computazionale, ma anche elevata complessità.

Copyright N.A. Borghese Università di Milano 26/03/2003

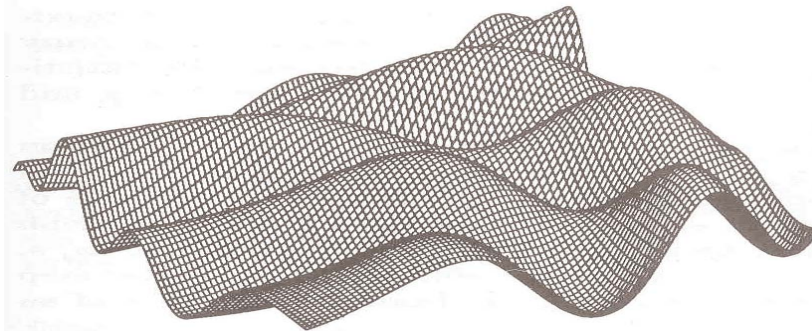
<http://homes.dsi.unimi.it/~borghese>



Problemi nell'apprendimento supervisionato



- Nota: $W_{iniziale}$ è generalmente casuale e può condizionare la convergenza degli algoritmi iterativi.
- I problemi di convergenza sono legati all'esistenza di minimi locali del funzionale J



Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Hebbian learning rule (1949)



...”When the axon of a cell A is near enough to excite a cell B and repeatedly or persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells such that A's efficiency, as one of the cells firing B, is increased....” ($\Delta w = f(y,x)$).

La forza di una sinapsi aumenta con l'utilizzo => Memoria?

Memoria a breve termine. Circuiti elettrici.

Memoria a lungo termine. Modificazioni chimiche.

In termini biologici si chiama **potenziamento**. LTP.

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Addestramento del perceptrone (learning rule)



Le uscite sono tutte indipendenti.

$$y_i = \text{sgn}(h_i) = \text{sgn}(\mathbf{w}_i \cdot \mathbf{u})$$

Soluzione incrementale:

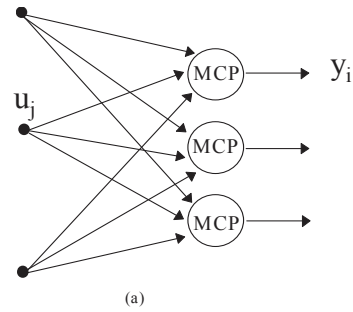
$$\mathbf{w}_{ij}^{\text{nuovo}} = \mathbf{w}_{ij}^{\text{vecchio}} + \Delta \mathbf{w}_{ij}$$

$$\Delta \mathbf{w}_{ij} = \eta (y_i^D - y_i) \mathbf{u}_j$$

oppure

$$\eta (1 - y_i^D y_i) y_i^D \mathbf{u}_j$$

Updating sse $y_i^D y_i < 0$



Hebbian learning

$$-1 < y_i < +1$$

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Addestramento del perceptrone con unita' a scalino



Le uscite sono tutte indipendenti.

$$y_i = \text{sgn}(h_i) = \text{sgn}(\mathbf{w}_i \cdot \mathbf{u})$$

Soluzione incrementale:

$$\mathbf{w}_{ij}^{\text{nuovo}} = \mathbf{w}_{ij}^{\text{vecchio}} + \Delta \mathbf{w}_{ij}$$

Margine di sicurezza n_k .

$$y_i^D y_i > 0 \implies y_i^D h_i > n_k \quad y_i^D \sum_j \mathbf{w}_{ij} \mathbf{u}_j > n_k$$

Rosenblatt perceptron learning rule:

Attivazione del neurone i

$$\Delta \mathbf{w}_{ij} = \eta \Theta(n_k - y_i^D h_i) y_i^D \mathbf{u}_j$$

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Analisi della funzione di addestramento (I)

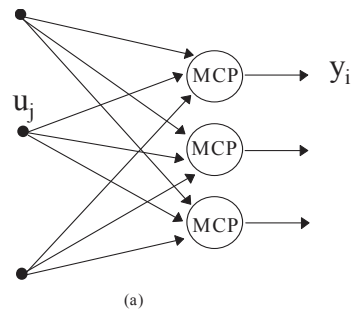


$$\Delta w_{ij} = \eta \Theta(n_k - y_i^D h_i) y_i^D u_j$$

positivo o nullo

n_k margine di errore

$$\Delta w_{ij} \neq 0 \iff \eta \Theta(1 - y_i^D h_i) y_i^D u_j$$



Iniziamo ad analizzare il termine $y_i^D h_i$ ponendo inizialmente $n_k=1$

Casi in cui y_i^D e h_i sono concordi (funzionamento corretto).

- 1) $h_i < 0$ $y_i = -1$ & $y_i^D = -1$: $y_i^D h_i > 0$
- 2) $h_i \geq 0$ $y_i = 1$ & $y_i^D = 1$: $y_i^D h_i \geq 0$



Analisi della funzione di addestramento (II)

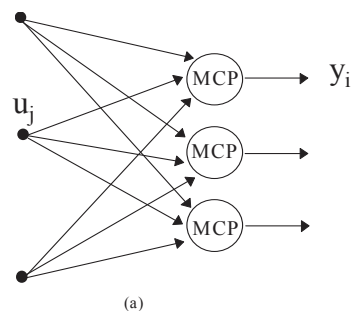


$$\Delta w_{ij} = \eta \Theta(n_k - y_i^D h_i) y_i^D u_j$$

negativo o nullo

$n_k = 0$

$$\Delta w_{ij} = 0 \iff \Theta(n_k - y_i^D h_i)$$



Iniziamo ad analizzare il termine $y_i^D h_i$ ponendo inizialmente $n_k=1$

Casi in cui y_i^D e h_i sono discordi (funzionamento scorretto).

- 1) $h_i < 0$ $y_i = -1$ & $y_i^D = 1$: $y_i^D h_i < 0$
- 2) $h_i \geq 0$ $y_i = 1$ & $y_i^D = -1$: $y_i^D h_i < 0$



Lo spirito Hebbiano



$$\Delta w_{ij} = \eta \Theta(n_k - y_i^D h_i) y_i^D u_j$$

$\Theta(\bullet)$ decide solo se la correzione deve essere effettuata (0,1).

Il segno dipende dal prodotto: $y_i^D u_j$

Il prodotto è positivo quando ingresso e uscita sono concordi. In questo caso tende a fare aumentare w_{ij} .

Il prodotto è negativo quando ingresso e uscita sono discordi. In questo caso tende a fare diventare più negativo w_{ij} .

L'algoritmo di apprendimento converge in un numero finito di passi (se la soluzione esiste!).



Apprendimento robusto (ruolo di n_k)



$$\Delta w_{ij} = \eta \Theta(n_k - y_i^D h_i) y_i^D u_j$$

Se y_i^D e h_i sono discordi, l'argomento è già in grado di attivare la funzione $\Theta(\bullet)$ anche se $n_k = 0$.

Se y_i^D e h_i sono concordi, perché $\Theta(\bullet)$ dia un valore = 1, occorre che: $|h_i| < n_k$.

n_k prende il nome di **margine di sicurezza** o **margine di errore** e stabilizza la rete rendendola più robusta al rumore sugli ingressi.

Nel caso di input e pesi binari, n_k conta di quanto $w_k u_k > w_j u_j$.



La pratica dell'apprendimento supervisionato



Fino a quando l'apprendimento non è stato completato:

1. Presentazione di un pattern di input / output.
2. Calcolo dell'output della rete con il pattern corrente.
3. Calcolo dell'incremento dei pesi.

Aggiornamento dei pesi.

Aggiornamento dei pesi:

- Per trial (ogni pattern)
- Per epoca (ogni insieme di pattern).



Ruolo di η – learning rate

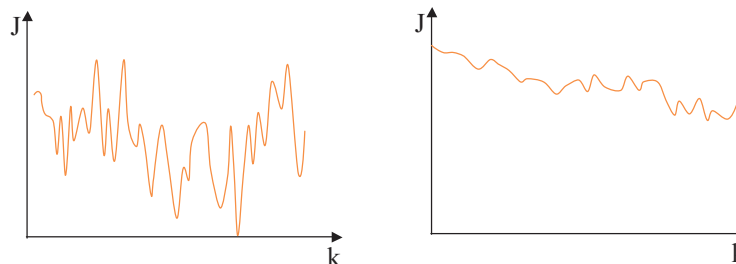


$$\Delta w_{ij} = \eta \Theta(n_k - y_i^D h_i) y_i^D u_j$$

Calmiera il Δw_{ij} per evitare che :

- Un peso sia specifico di un'unità ingresso-uscita.
- Oscillazioni durante l'apprendimento senza convergenza.

η può variare durante l'addestramento.





Perceptrone con unità di attivazione continue



Possiamo derivare una regola di apprendimento di spirito Hebbiano per una qualsiasi funzione di attivazione continua

$$y = g\left(\sum_{j=1} (w_{ij} u_j - \mu_i)\right) = g\left(\sum_{j=0} (w_{ij} u_j)\right)$$

Si tratta di un problema di minimizzazione di una cifra di merito, J , sullo spazio di parametri W :

$$J = \underbrace{\|y^D - g(W^{nuovo} U)\|}_{\text{Errore}} \leq \|y^D - g(W^{vecchio} U)\|$$

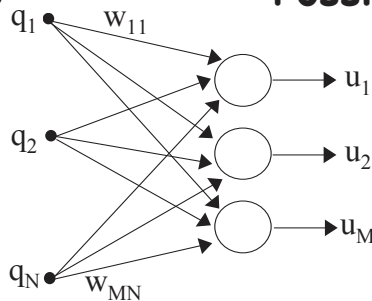
$$J = E(w) = \frac{1}{2} \sum_p \left[\sum_i (y_{ip}^D - y_{ip})^2 \right] = \frac{1}{2} \sum_i \left(y_{ip}^D - g\left(\sum_j w_{ij} u_{jp}\right) \right)^2$$

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Possibili soluzioni



$$\bar{w} = \begin{vmatrix} w_{11} & w_{12} & \dots & w_{1N} \\ w_{21} & w_{22} & \dots & \dots \\ \dots & \dots & \dots & \dots \\ w_{M1} & w_{M2} & \dots & w_{MN} \end{vmatrix} M \times N$$

Dimensioni dello spazio di ricerca

Possibili soluzioni del problema di minimizzazione:

- Gradiente e sue modificazioni
- Algoritmi Newtoniani e quasi Newtoniani
- Algoritmi genetici
- Altri algoritmi di ricerca operativa/calcolo numerico

Trovare ΔW tale che l'insieme dei pesi $W^{nuovo} = W^{vecchio} + \Delta W$ dia luogo a un'uscita più vicina all'uscita desiderata di $W^{vecchio}$

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Unità di attivazione lineari



$$y = g\left(\sum_{j=1} (w_{ij} u_j - \mu_i)\right) = g\left(\sum_{j=0} (w_{ij} u_j)\right)$$

Caso lineare ($g = 1$):

$$y_i = \sum_{j=1} (w_{ij} u_j - \mu_i) = \sum_{j=0} (w_{ij} u_j) \implies \mathbf{Y} = \mathbf{W} \mathbf{U}$$

Soluzione di un sistema lineare nei pesi!!

Condizione di risolubilità: \mathbf{U} di rango massimo \rightarrow
 $\{w\}$ sono linearmente indipendenti.



Unità lineari, soluzione iterativa



$$J = E(\mathbf{w}) = \frac{1}{2} \sum_p \left[\sum_i (y_{ip}^D - y_{ip})^2 = \frac{1}{2} \sum_i \left(y_{ip}^D - \left(\sum_j w_{ij} u_{jp} \right) \right)^2 \right]$$

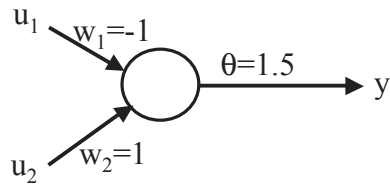
$$\Delta w_{ijp} = -\eta \frac{\partial}{\partial w_{ij}} \frac{1}{2} \sum_i \left(y_i^D - \left(\sum_j w_{ij} u_j \right) \right)^2$$

$$\Delta w_{ijp} = +\eta \sum_i \left(y_i^D - \left(\sum_j w_{ij} u_j \right) \right) u_j = +\eta (y_i^D - y_i) u_j$$

δ rule (1960)



Esempio di delta rule - I



$$U = [-1, 1] \quad y^D = -1 \\ \eta = 0.2$$

u_1	u_2	y^D
-1	-1	-1
-1	1	-1
1	-1	-1
1	1	1

$$y = \sum_{j=1} (w_j u_j - \theta) = \sum_{j=0} (w_j u_j) = (-1)(-1) + (1)(1) - 1.5 = 0.5 \gg -1$$

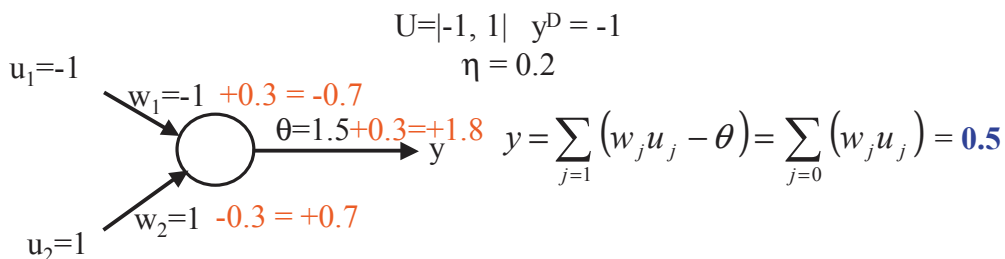
$$u_0 = 1 \quad w_0 = -\theta$$

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Esempio di delta rule - II



$$U = [-1, 1] \quad y^D = -1 \\ \eta = 0.2$$

$$y = \sum_{j=1} (w_j u_j - \theta) = \sum_{j=0} (w_j u_j) = 0.5$$

$$\Delta w_{ij} = +\eta (y_i^D - y_i) u_j$$

$$-\Delta \theta = \Delta w_0 = \eta (y_i^D - y_i) u_0 = \eta (-1 - 0.5)(1) = -0.30$$

$$\Delta w_1 = \eta (y_i^D - y_i) u_1 = \eta (-1 - 0.5)(-1) = +0.30$$

$$\Delta w_2 = \eta (y_i^D - y_i) u_2 = \eta (-1 - 0.5)(1) = -0.30$$

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Esempio di delta rule - III



$U = \{-1, 1\} \quad y^D = -1$
 $\eta = 0.2$

$u_1 = -1$
 $w_1 = -0.7 \quad +0.12 = -0.58$
 $u_2 = 1$
 $w_2 = 0.7 \quad -0.12 = +0.58$

$\theta = 1.8 + 0.12 = +1.92$

$y = \sum_{j=1} (w_j u_j - \theta) = \sum_{j=0} (w_j u_j) = -0.4$
 -0.76

$\Delta w_{ij} = +\eta (y_i^D - y_i) \mu_j$

$-\Delta \theta = \Delta w_0 = \eta (y_i^D - y_i) \mu_0 = \eta (-1 - (-0.4))(1) = -0.12$
 $\Delta w_1 = \eta (y_i^D - y_i) \mu_1 = \eta (-1 - (-0.4))(-1) = +0.12$
 $\Delta w_2 = \eta (y_i^D - y_i) \mu_2 = \eta (-1 - (-0.4))(1) = -0.12$

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Esempio di delta rule - Cattiva scelta di η



$U = \{-1, 1\} \quad y^D = -1$
 $\eta = 0.8$

$u_1 = -1$
 $w_1 = -1 \quad +1.2 = +0.2$
 $u_2 = 1$
 $w_2 = 1 \quad -1.2 = -0.2$

$\theta = 1.5 + 1.2 = +2.7$

$y = \sum_{j=1} (w_j u_j - \theta) = \sum_{j=0} (w_j u_j) = 0.5$
 -2.3

$\Delta w_{ij} = +\eta (y_i^D - y_i) \mu_j$

$-\Delta \theta = \Delta w_0 = \eta (y_i^D - y_i) \mu_0 = \eta (-1 - 0.5)(1) = -1.2$
 $\Delta w_1 = \eta (y_i^D - y_i) \mu_1 = \eta (-1 - 0.5)(-1) = +1.2$
 $\Delta w_2 = \eta (y_i^D - y_i) \mu_2 = \eta (-1 - 0.5)(1) = -1.2$

Copyright N.A. Borghese Università di Milano 26/03/2003

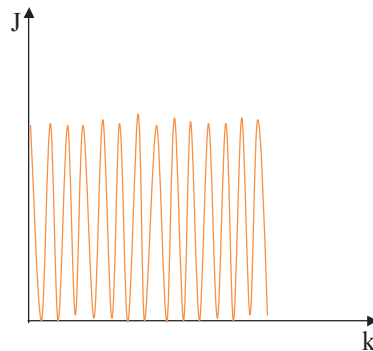
<http://homes.dsi.unimi.it/~borghese>



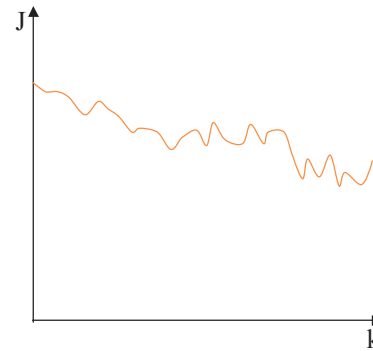
Variazioni della funzione costo per diversi valori di η



η elevato



η piccolo

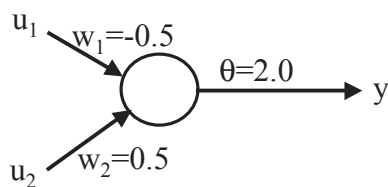


Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Esempio di specializzazione su un pattern



u_1	u_2	y^D
-1	-1	-1
-1	1	-1
1	-1	-1
1	1	1

a
b
c
d

a $y = \sum (w_j u_j - \theta) = \sum (w_j u_j) = (-0.5)(-1) + (-0.5)(1) - 2.0 = -2$

b $y = \sum_{j=1}^1 (w_j u_j - \theta) = \sum_{j=0}^0 (w_j u_j) = (-0.5)(-1) + (0.5)(1) - 2.0 = -1$

c $y = \sum (w_j u_j - \theta) = \sum (w_j u_j) = (-0.5)(1) + (0.5)(-1) - 2.0 = -3$

d $y = \sum_{j=1}^1 (w_j u_j - \theta) = \sum_{j=0}^0 (w_j u_j) = (-0.5)(1) + (0.5)(1) - 2.0 = -2$

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



Unità non-lineari, soluzione iterativa



$$J = E(\mathbf{w}) = \frac{1}{2} \sum_p \left[\sum_i (y_{ip}^D - y_{ip})^2 \right] = \frac{1}{2} \sum_p \left[\sum_i \left(y_{ip}^D - g \left(\sum_j w_{ij} u_{jp} \right) \right)^2 \right]$$

$$\Delta w_{ijp} = -\eta \frac{\partial}{\partial w_{ij}} \frac{1}{2} \sum_i \left(y_i^D - g \left(\sum_j w_{ij} u_j \right) \right)^2 =$$

$$\eta \sum_i \left(y_i^D - g \left(\sum_j w_{ij} u_j \right) \right) g' \left(\sum_j w_{ij} u_j \right) u_j = +\eta \underbrace{\left(y_i^D - y_i \right) g' \left(\sum_j w_{ij} u_j \right)}_{\delta \text{ rule}} u_j$$

Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>



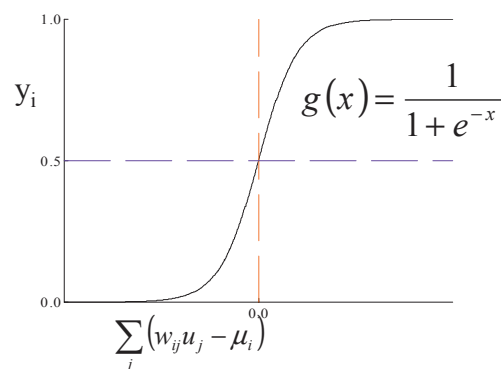
Perceptrone con unità di attivazione logistiche



$$g'(x) = g(x) \cdot (1 - g(x))$$

$$y_i = g \left(\sum_j (w_{ij} u_j - \mu_i) \right)$$

$$\begin{aligned} g'(x) &= \frac{e^{-x}}{(1 + e^{-x})^2} = \\ &= \frac{1}{1 + e^{-x}} \left(1 - \frac{1}{1 + e^{-x}} \right) \end{aligned}$$



Copyright N.A. Borghese Università di Milano 26/03/2003

<http://homes.dsi.unimi.it/~borghese>

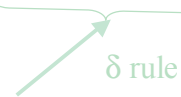


Update dei pesi per funzione logistica



$$J = E(\mathbf{w}) = \frac{1}{2} \sum_p \left[\sum_i (y_{ip}^D - y_{ip})^2 = \frac{1}{2} \sum_i \left(y_{ip}^D - g\left(\sum_j w_{ij} u_{jp}\right) \right)^2 \right]$$

$$\Delta w_{ijp} = +\eta \sum_i (y_i^D - g(\cdot)) g'(\cdot) u_j = +\eta \underbrace{(y_i^D - y_i)}_{\delta \text{ rule}} y_i (1 - y_i) u_j$$

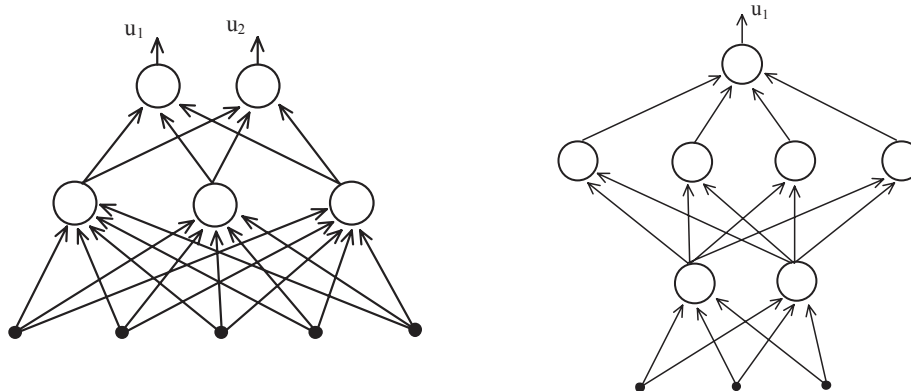


Esistenza della soluzione: indipendenza lineare dei pattern.

La funzione costo non è quadratica → Minimi locali.



Perceptrone a più strati



Algoritmi di *apprendimento più sofisticati: back-propagation.*

Collegamento con la teoria statistica dell'apprendimento

→ *Corso di reti neurali.*



Riassunto - Apprendimento



Algoritmi iterativi per adattare il valore dei parametri (pesi).

Definizione di una funzione costo che misura la differenza tra valore fornito e quello desiderato.

Algoritmo (gradiente) che consente di aggiornare i pesi in modo da minimizzare la funzione costo.

Training per pattern (specializzazione) o per epoche.



Problemi



Qual è il problema principale dell'apprendimento supervisionato?

L'uscita delle funzioni logistiche è compresa tra 0 e 1. Come si possono approssimare funzioni con un range più ampio?





Problemi



Quando si termina l'algoritmo di apprendimento?

Bootstrap – Vengono estratti pattern con ripetizioni.

Cross-Validation - Errore sull'insieme di training =
Errore sull'insieme di test.

Utilizzare lo “structural risk” invece dell’“empirical risk”.

Si vuole evitare che la rete si specializzi troppo sui pattern di training e non sia in grado di interpolare.

Applicazioni alla biologia: esperimento di Zipser e Andersen sui neuroni del collicolo superiore.