

COMPENDIO DI RETI NEURALI

- A.A 2018/2019

A cura di Cesare Iurlaro

<https://linkedin.com/in/cesare-iurlaro>

cesareiurlaro@gmail.com

1 – Percettrone

- Funzione di attivazione e architettura
- *Algoritmo di apprendimento*
 - o **Delta rule**
 - o **Teorema di convergenza: enunciato e dimostrazione**
- Limiti del percertrone

2 – Adaline

- Funzione di attivazione e architettura
- *Algoritmo di apprendimento LMS*
 - o **Filtering adattivo**
 - o **Funzione costo**
 - o **Problema unconstraint optimization (definizione + condizione di ottimalità)**

 - **Approccio local iterative descent**

 - **Steepest Descent (Discesa del gradiente)**

 - o **Stima del gradiente in LMS**

3 – MLP

- Funzione di attivazione e architettura
- *Algoritmo di apprendimento Back-Propagation*
 - o **Forward/Backward phase**
 - o **Funzioni costo: istantaneo e medio**
 - o **Stima del gradiente in Back-Propagation**
 - **Gradiente locale output/hidden**
- Metodi di apprendimento (epoche/pattern)
- Problemi di η e soluzioni
- Condizioni di terminazione (variazione dell'errore o generalizzazione)
- Topologia della rete (pruning/aggiunta iterativa)
- Rete come approssimatore di funzioni (booleane, continue a tratti, continue)
- **Teorema di esistenza dell'approssimazione universale**

4 – RBF

- Funzione di attivazione e architettura
 - o Parametri
 - o Approccio ibrido
- Teorema di Cover
- Problema dell'interpolazione
- Teorema di Micchelli
 - o Funzioni per cui è valido il teorema
- *Algoritmo di apprendimento (determinazione dei parametri)*
 - o **Metodo della pseudoinversa ottimizzato per calcolo parallelo e distribuito**
 - o **Hybrid Learning Process**

5 – ELM

- Funzione di attivazione e architettura
- Teoremi di ELM
- *Algoritmo di apprendimento*
 - o **Pseudoinversa**
 - o **Pseudoinversa di SVD per gestire singolarità, evitando l'inversione matriciale**

6 – CNN

- Introduzione
 - o Studio delle reti neurali
 - o Arrivo di AlexNet
 - o Fattori del suo successo
- Funzione di attivazione e architettura
 - o **Convoluzione e livelli convoluzionali**
 - o **Pooling**
 - o **Relu e vantaggi**
- Scelte in architetture recenti
 - o Soft-max layer
 - o Cross-entropy loss function
- *Apprendimento*
 - o **Transfer learning**
 - o **Features re-using**

7 – SOM

- Architettura
- Processi di una SOM
 - o **Competizione** (discriminazione del *BMU*)
 - o **Cooperazione** (selezione del *vicinato topologico*)
 - o **Adattamento sinaptico** (aggiornamento vettore pesi)
 - **Fasi del processo adattivo (ordinamento e convergenza)**
- *Algoritmo di apprendimento*
 - o **Decadimento esponenziale iterativo di η e σ**

8 – Hopfield & Boltzmann

- Definizione di energy based model
 - o **Funzione di energia globale**
 - o **Energy gap**
 - o **Memoria stabile e memoria fondamentale**
- Pattern completion e algoritmo non deterministico di information withdrawal
- Teorema di convergenza
- Principio di Hebb
- *Algoritmo di apprendimento*
 - o **Limitazioni**
 - o **Boltzmann: Simulated annealing**
- *RBM*
 - o **Architettura**
 - o **Algoritmo di apprendimento**
- *Esempi di calcolo*

0 – Introduzione alle Reti Neurali (Neural Networks)

Elemento computazionale massicciamente distribuito e parallelo che memorizza in un vettore pesi conoscenza sperimentale.

Allenamento: processo di miglioramento delle prestazioni, realizzato sottoponendo la rete ad esempi (conoscenza sperimentale).

Benefici:

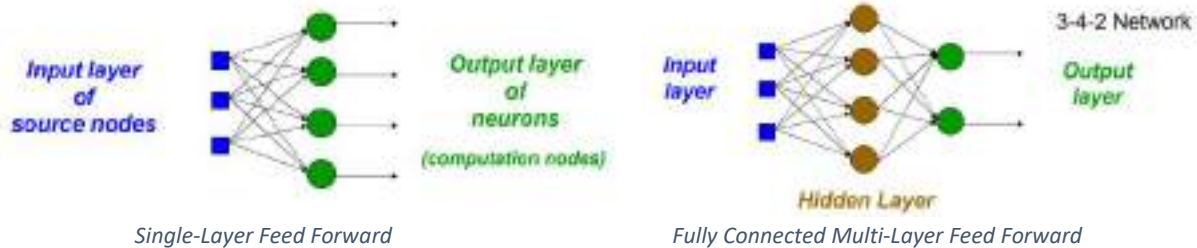
- Non-linearity
- Input-Output mapping realizzato attraverso l'apprendimento.
- Fault tolerance: le prestazioni degradano gradualmente.
- Incremental Training: è possibile continuare l'apprendimento se sono disponibili nuovi dati.

Componenti di una *Neural Network*:

- Architettura: insieme di neuroni collegamenti tramite link. A ciascun link è associato un peso.
- Neurone: unità computazionale dell'*architettura*.
- Algoritmo di apprendimento: algoritmo usato per allenare la *Neural Network*.

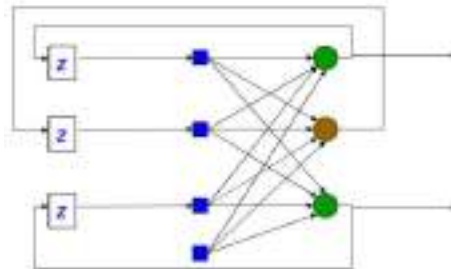
Classi di architettura:

- Feed Forward (Single-Layer o Multi-Layer): i neuroni sono organizzati a livelli in un grafo aciclico



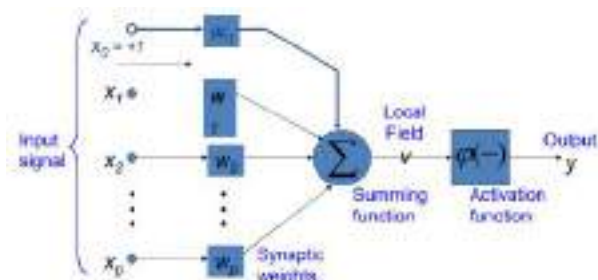
La presenza di *Hidden Layer* permette di estrarre dinamiche di più alto ordine proiettando i dati in spazi differenti. Ciò assume particolare importanza all'aumentare dei nodi al *livello di input*.

- Recurrent: architettura caratterizzata da **unit delay operators**, componenti che controllano dinamicamente il sistema. Rendono possibili dinamiche temporali.



Recurrent Neural Network

Architettura del neurone:



- **Link**: caratterizzano l'input del neurone associato a un peso w_i .
- **Combinazione lineare u_j** : funzione che computa la somma degli input pesata in funzione dei pesi.

$$u_j = \sum_{i=1}^p w_{ji} \cdot x_i$$

- **Bias $w_{j0} = b_j$** : intercetta che applica una trasformazione alla combinazione lineare u_j : $v_j = u_j - b_j$.

- **Funzione di attivazione**: funzione che limita l'intervallo dei valori possibili che l'output può assumere.

$$y_j = \varphi(u_j - b_j) = \varphi(v_j)$$

La scelta della funzione φ (*step*, *ramp*, *sigmoid* o *gaussian*) caratterizza il modello del neurone.

1 – Percettrone

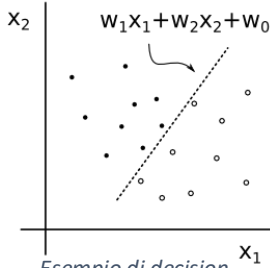
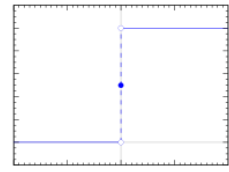
Il più semplice *modello* di Rete Neurale; classifica esempi in uno spazio linearmente separabile.

Induced local field $v_j = \mathbf{w}^T \mathbf{x}$: combinazione lineare del **vettore pesi** \mathbf{w} e delle p componenti di \mathbf{x} .
Il vettore \mathbf{x} è detto *esempio, istanza o pattern input*.

Il percettrone è caratterizzato dalla **funzione di attivazione a gradino**: $\varphi(v_j) = \begin{cases} 1 & v_j > 0 \\ -1 & v_j < 0 \end{cases}$

Dati il vettore pesi \mathbf{w} , è possibile visualizzare il **decision boundary** (ovvero l'iperpiano che divide le istanze in funzione delle classi di appartenenza) tracciato sullo spazio delle istanze.

Funzione φ a gradino



Esempio di decision boundary in due dimensioni

Apprendimento del percettrone – Supervisionato

Sia $\mathbf{w}(t)$ il vettore pesi \mathbf{w} all'iterazione t ed η il *learning rate* (coefficiente di apprendimento).

Sia inoltre la seguente **regola di aggiornamento**: $\mathbf{w}(k+1) = \mathbf{w}(k) + \eta \mathbf{x}(k) d(k)$

$$\Rightarrow \Delta \mathbf{w} = \mathbf{w}(k+1) - \mathbf{w}(k) = \eta \mathbf{x}(k) d(k)$$

L'algoritmo di apprendimento del percettrone modifica iterativamente \mathbf{w} nel seguente modo:

- 1) Inizializza \mathbf{w} con valori casuali
- 2) While (la presente iterazione n presenta esempi erroneamente classificati) {
 1. Seleziona un esempio classificato male
 2. Aggiorna i pesi con la *regola di aggiornamento*

Osservazione: L'algoritmo converge quando la condizione di terminazione viene rispettata, ovvero quando \mathbf{w} è tale per cui il decision boundary classifichi correttamente **tutti** gli esempi.

Teorema di convergenza: Siano le classi C_1, C_2 linearmente separabili; allora l'algoritmo di apprendimento del percettrone converge in un numero finito di iterazioni, ovvero esiste un iperpiano che separi correttamente C_1 e C_2 .

Dimostrazione:

Assumiamo $\eta = 1$ e $\mathbf{w}(0) = \mathbf{0}$.

Sia C_1 l'insieme delle istanze per cui $d_{C_1} = 1$ e C_2 l'insieme di istanze per cui $d_{C_2} = -1$.

$d(\mathbf{x})$: *desired output dell'istanza* \mathbf{x} .

Sia C l'insieme ottenuto dall'unione $C_1 \cup C_2$.

- 1) **Passo preliminare di correzione:** $\forall \mathbf{x}(i) \in C_2$ applichiamo la trasformazione $\mathbf{x}(i) = -\mathbf{x}(i)$, cioè tale che $\forall n, d(n) = 1$
- 2) **Definizione del problema di ricerca:** trovare un vettore pesi \mathbf{w}_* tale che $\forall \mathbf{x}(i) \in C, \mathbf{w}_*^T \mathbf{x}(i) > 0$, cioè tale che $\forall \mathbf{x}(i), y(i) = 1$
- 3) Sia k la k -esima iterazione. Poiché:
 - a. $\forall n, d(n) = 1$
 - b. $\eta = 1$

Allora per la regola di aggiornamento $\mathbf{w}(k) + \eta \mathbf{x}(k) d(k)$, vale che:

$$\forall k : \mathbf{w}(k+1) = \mathbf{w}(k) + \mathbf{x}(k) = \mathbf{x}(1) + \mathbf{x}(2) + \dots + \mathbf{x}(k)$$

A) Determinazione del lower bound; sia $\alpha = \min \|\mathbf{x}(i)\|$ l'istanza con il modulo minore.

$$\mathbf{w}(k+1) = \mathbf{x}(1) + \mathbf{x}(2) + \dots + \mathbf{x}(k)$$

$$\mathbf{w}_*^T \cdot \mathbf{w}(k+1) = \mathbf{w}_*^T \cdot [\mathbf{x}(1) + \mathbf{x}(2) + \dots + \mathbf{x}(k)]$$

$$\mathbf{w}_*^T \cdot \mathbf{w}(k+1) \geq k\alpha$$

$$\|\mathbf{w}_*^T \cdot \mathbf{w}(k+1)\|^2 \geq (k\alpha)^2$$

$$\text{Disuguaglianza di Cauchy-Schwartz: } |(u, v)| \leq \|u\| \cdot \|v\| \Rightarrow \|\mathbf{w}_*^T\|^2 \cdot \|\mathbf{w}(k+1)\|^2 \geq \|\mathbf{w}_*^T \cdot \mathbf{w}(k+1)\|^2 \geq (k\alpha)^2$$

$$\|\mathbf{w}(k+1)\|^2 \geq \left(\frac{k\alpha}{\|\mathbf{w}_*^T\|} \right)^2$$

B) Determinazione dell'upper bound; sia $\beta = \max \|\mathbf{x}(i)\|^2$ l'istanza con il modulo maggiore.

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \mathbf{x}(k)$$

$$\|\mathbf{w}(k+1)\|^2 = \|\mathbf{w}(k) + \mathbf{x}(k)\|^2$$

$$\|\mathbf{w}(k+1)\|^2 = \|\mathbf{w}(k)\|^2 + \|\mathbf{x}(k)\|^2 + 2 \cdot \mathbf{w}^T(k) \cdot \mathbf{x}(k)$$

$$\mathbf{w}(k+1) \Rightarrow \mathbf{x}(k) \text{ non classificato correttamente} \Rightarrow 2 \cdot \mathbf{w}^T(k) \cdot \mathbf{x}(k) \leq 0$$

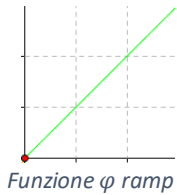
$$\Rightarrow \|\mathbf{w}(k+1)\|^2 \leq \|\mathbf{w}(k)\|^2 + \|\mathbf{x}(k)\|^2 = \sum_{i=0}^k \|\mathbf{x}(i)\|^2 \Rightarrow \|\mathbf{w}(k+1)\|^2 \leq \sum_{i=0}^k \|\mathbf{x}(i)\|^2$$

$$\|\mathbf{w}(k+1)\|^2 \leq k\beta$$

- 4) Da **A** e **B** segue che: $\left(\frac{k\alpha}{\|\mathbf{w}_*^T\|} \right)^2 \leq \|\mathbf{w}(k+1)\|^2 \leq k\beta$, ovvero che l'algoritmo termina in un numero finito di iterazioni.

Limite del percettrone: non risolve problemi non linearmente separabili (es: *problema dello XOR*)

2 – Adaline: Adaptive Linear Element



Modello di Rete Neurale caratterizzato dalla funzione di attivazione lineare ramp $y = \phi(v)$.

Utilizza l'algoritmo di apprendimento supervisionato con filtering adattivo LMS (Least Mean Square).

Filtering process: computazione dei segnali di *output* $y(i)$ e di *errore* $e(i)$, dato un vettore *input* x_i .

Adaptive process: aggiustamento automatico del vettore pesi w_i in funzione dell'errore $e(i)$.

Problema unconstraint-optimization: calcolare $w^* = \operatorname{argmin}_w E(w)$. $E(w)$: *funzione costo* che misura la bontà di un vettore w

Condizione necessaria per l'ottimalità: $\nabla E(w^*) = 0$ (notare che 0 è un vettore di zeri)

○ ∇ è l'**operatore gradiente**

○ $\nabla E(w) = \left[\frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_m} \right]$ è il **gradiente** (o **vettore gradiente**) della funzione costo

Gradiente di una funzione: vettore le cui componenti sono le derivate parziali della funzione rispetto agli assi di riferimento.

- La **direzione** del gradiente identifica la *direzione più ripida* della funzione.

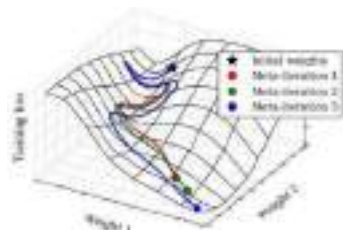
Ciascuna componente j del gradiente definisce due informazioni:

1. La **magnitudine** con cui ciascun peso w_j influenza l'errore $E(w)$.
2. Il **segno** necessario a w_j per ridurre l'errore (se w_j è più grande o piccolo di come dovrebbe essere).

Local iterative descent: approccio per problemi di *unconstraint-optimization*.

Consiste nella *riduzione iterativa* della *funzione costo*. Ad ogni k -esima iterazione, vale che: $E(w(k+1)) < E(w(k))$.

- **Metodo Steepest Descent (discesa del gradiente):** tecnica di ottimizzazione con approccio di *local iterative descent* che iterativamente modifica il vettore pesi $w(k)$ al fine di massimizzare la riduzione dell'errore in esso presente.



Esempio di discesa del gradiente

La *regola di aggiornamento* di questo metodo prende il nome di **steepest descent delta rule**:

$$w(k+1) = w(k) - \eta(\nabla E(w))$$

$$\Rightarrow \Delta w = w(k+1) - w(k) = -\eta(\nabla E(w))$$

Ossia, ad ogni iterazione, w viene aggiornato in funzione del gradiente della funzione costo (stessa direzione e verso opposto) ponderato su $\eta \in \mathbb{R}$, che è il tasso di apprendimento.

➤ **Algoritmo LMS (1960):** algoritmo che effettua *filtering adattivo lineare*. Usa il *metodo steepest descent* per **minimizzare** il **valore** della funzione costo quadratico istantaneo $E(w) = \frac{1}{2}e^2(k)$, attraverso la **stima** del vettore gradiente *istantaneo* $\hat{g}(k)$:

Notare che l'algoritmo minimizza l'errore non attraverso il calcolo del gradiente effettivo, ma di una sua stima! Pertanto anche il vettore pesi calcolato è una stima, e viene indicato come \hat{w} .

$$E(\hat{w}) = \frac{1}{2}e^2(k) \Rightarrow \frac{\partial E(\hat{w})}{\partial e(k)} = \frac{\partial \left(\frac{1}{2}e^2(k) \right)}{\partial e(k)} = e(k)$$

$$e(k) = d(k) - \hat{w}^T(k)x(k) \Rightarrow \frac{\partial e(k)}{\partial \hat{w}(k)} = -x(k)$$

$$\hat{g}(k) = \frac{\partial E(\hat{w})}{\partial \hat{w}} = \frac{\partial E(\hat{w})}{\partial e(k)} \cdot \frac{\partial e(k)}{\partial \hat{w}} = -x(k)e(k)$$

Dalla *steepest descent delta rule* segue la **LMS delta rule**:

$$\hat{w}(k+1) = \hat{w}(k) + \eta x(k) e(k)$$

$$\Rightarrow \Delta \hat{w} = \hat{w}(k+1) - \hat{w}(k) = \eta x(k) e(k)$$

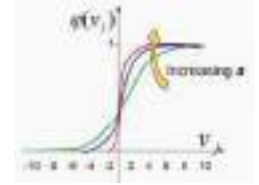
Algoritmo LMS:

- 1) Inizializza casualmente $w(k)$.
- 2) While (E_{tot} insoddisfacente $\wedge k < \max_iterations$) {
 1. Itera su un esempio $k : (x^{(k)}, d^{(k)})$
 2. Computa l'**errore istantaneo** k -esimo : $e^{(k)} = d(k) - x^T(k) \hat{w}(k)$
 3. Aggiorna i pesi con la **LMS Delta rule**: $w(k+1) = w(k) + \eta x(k) e^{(k)}$

3 – MLP: Multilayer Perceptron

La struttura di **MLP** è definita tramite una sequenza di *layer* di *perceptroni* e gode delle seguenti caratteristiche:

- Il modello di ciascun neurone nella rete include una funzione di attivazione *non lineare e differenziabile*, tipicamente la funzione di attivazione *sigmoide* $\varphi(v) = \frac{1}{1+e^{-av}}$, dove $v = \mathbf{w}_{ji} \boldsymbol{\varphi}$.
- **Uno o più hidden layer**: nodi che non sono né di input né di output.
- Esibisce un alto livello di connettività.



Funzione φ sigmoide

Gli *hidden layer* effettuano una **trasformazione non lineare** dei dati in input in uno spazio chiamato *spazio delle feature*.

Rete neurale feed-forward (FFNN): modello di rete neurale dove le connessioni tra le unità non formano cicli.

Tipicamente l'algoritmo di *apprendimento supervisionato* in una rete di questo tipo è l'**algoritmo di back-propagation (1985)**, che in particolari degenera nell'*algoritmo LMS*. Tale algoritmo consiste in due iterazioni:

- 1) **Forward phase**: fissati i pesi della rete, il segnale di input si propaga attraverso la rete fino all'output.
- 2) **Backward phase**: l'output della rete viene comparato con la risposta attesa e ne viene computato l'errore. Tale errore viene propagato all'indietro per la rete e ne modifica i pesi.

- $\mathcal{T} = \{\mathbf{x}(k), \mathbf{d}(k)\}_{k=1}^N$ il **training set** di N istanze, etichettato.
- $y_j(k)$ l'**output** del **j-esimo neurone nell'output layer** allo stimolo dell'istanza $\mathbf{x}(k)$.
- $e_j(k) = d_j(k) - y_j(k)$ l'**errore** del **j-esimo neurone nell'output layer** allo stimolo dell'istanza $\mathbf{x}(k)$.
- $E(k) = \frac{1}{2} \sum_j e_j^2(k)$ l'**errore totale istantaneo** dell'intero *output layer* allo stimolo dell'istanza $\mathbf{x}(k)$.
- $E_{av}(N) = \frac{1}{N} \sum_{k=1}^N E(k)$ l'**errore medio sul training set** (o **rischio empirico**)

Osservazione: $E(k)$ è la somma, per ogni *neurone output*, della funzione costo che era presente in *Adeline*.

Stima del gradiente nell'algoritmo di back-propagation:

$$\begin{aligned} \frac{\partial E(\mathbf{w})}{\partial e_j(k)} &= e_j & \frac{\partial e_j(k)}{\partial y_j(k)} &= -1 \\ \frac{\partial y_j(k)}{\partial v_j(k)} &= \varphi'_j(v_j) & \frac{\partial v_j(k)}{\partial \mathbf{w}_{ji}(k)} &= y_i \end{aligned}$$

A seconda della tipologia di *layer* a cui appartengono i neuroni (*hidden* oppure *output*), il *gradiente locale* assume forme diverse. Questo è vero poiché agli *output layer* sono associate informazioni che non lo sono agli *hidden layer*: i *desired output* $d(n)$.

- **Gradiente computato nei neuroni output:**

$$\hat{g}_{out}(k) = \frac{\partial E}{\partial \mathbf{w}_{ji}} = \frac{\partial E}{\partial e_j} \cdot \frac{\partial e_j}{\partial y_j} \cdot \frac{\partial y_j}{\partial v_j} \cdot \frac{\partial v_j}{\partial \mathbf{w}_{ji}} = -e_j \varphi'_j(v_j) y_i$$

Dalla *steepest descent delta rule* segue la **Back-Propagation delta rule** (per i neuroni di *output*):

$$\hat{\mathbf{w}}_{ji}(k+1) = \hat{\mathbf{w}}_{ji}(k) + \eta e_j \varphi'_j(v_j) y_i$$

Il fattore $e_j \varphi'_j(v_j)$ prende il nome di **gradiente locale** (o **segnale d'errore**) del *j-esimo neurone di output*; si indica con δ_j :
 $\Rightarrow \Delta \mathbf{w}_{ji} = \mathbf{w}_{ji}(k+1) - \mathbf{w}_{ji}(k) = \eta \delta_j y_i$

- **Gradiente computato nei neuroni hidden:**

Q : insieme dei neuroni del livello successivo; $q \in Q$

$$\frac{\partial E}{\partial y_j} = \sum_q e_q \cdot \frac{\partial e_q}{\partial y_j} = \sum_q e_q \cdot \frac{\partial e_q}{\partial y_q} \cdot \frac{\partial y_q}{\partial v_j} \cdot \frac{\partial v_q}{\partial y_j} = \sum_q e_q \cdot [-1 \cdot \varphi'_q(v_q)] \cdot \hat{\mathbf{w}}_{qj} = -\sum_q e_q \varphi'_q(v_q) \hat{\mathbf{w}}_{qj} \stackrel{\text{def}}{=} -\sum_q \delta_q \hat{\mathbf{w}}_{qj}$$

$$\hat{g}_{hid}(k) = \frac{\partial E}{\partial \mathbf{w}_{ji}} = \frac{\partial E}{\partial y_j} \cdot \frac{\partial y_j}{\partial v_j} \cdot \frac{\partial v_j}{\partial \mathbf{w}_{ji}} = \left[-\sum_q \delta_q \hat{\mathbf{w}}_{qj} \right] \cdot \varphi'_j(v_j) \cdot y_i$$

Dalla *steepest descent delta rule* segue la **Back-Propagation delta rule** (per i neuroni *hidden*):

$$\hat{\mathbf{w}}_{ji}(k+1) = \hat{\mathbf{w}}_{ji}(k) + \eta \left(\sum_q \delta_q \hat{\mathbf{w}}_{qj} \right) \cdot \varphi'_j(v_j) \cdot y_i$$

Il fattore $\left(\sum_q \delta_q \hat{\mathbf{w}}_{qj} \right) \cdot \varphi'_j(v_j)$ prende il nome di **gradiente locale** del *j-esimo neurone hidden*; si indica con δ_j :

$$\Rightarrow \Delta \mathbf{w}_{ji} = \mathbf{w}_{ji}(k+1) - \mathbf{w}_{ji}(k) = \eta \delta_j y_i$$

Metodi apprendimento:

- Per epoche (batch): i pesi vengono modificati dopo la presentazione (*non ordinata*) di tutti gli N esempi nel *training set* \mathcal{T} . Un'intera presentazione prende il nome di **epoca di allenamento**. In altre parole, la *funzione costo* è definita da $E_{av}(N)$.
- Per pattern (online): i pesi vengono modificati dopo la presentazione *ordinata* di ciascun esempio nel *training set* \mathcal{T} . La *funzione costo* è definita da $E(k)$.

Problema: η piccolo implica apprendimenti lenti, η grande implica cambiamenti più consistenti ma potrebbe causare comportamenti instabili con oscillazioni del vettore pesi se > 1 .

Tecniche per aggirare il problema:

1. Introdurre un **termine momentum** $\alpha \Delta \mathbf{w}_{ji}(k-1)$ nella regola di aggiornamento dei pesi che prenda in considerazione gli aggiornamenti precedenti. Tale nuova regola prende il nome di **Delta rule generalizzata**:

$$\Delta \mathbf{w}_{ji}(k) = \alpha \Delta \mathbf{w}_{ji}(k-1) + \eta \delta_j \mathbf{y}_i ; \quad 0 \leq \alpha < 1$$

Il *termine momentum* ha l'effetto di accelerare l'apprendimento in presenza di discesa costante e di stabilizzarlo in presenza di oscillazioni.

2. **Aggiustamento adattivo di η** attraverso l'uso di euristiche:
 - a. *Euristica 1*: assegnare un η diverso a ciascun peso
 - b. *Euristica 2*: variare η da un'iterazione all'altra successiva

L'apprendimento continua finché la condizione di terminazione non è soddisfatta.

Possibili condizioni di terminazione:

- Variazioni dell'errore quadratico medio totale:
L'algoritmo converge se il tasso assoluto di cambiamento di $E_{av}(N)$ per epoca ricade nell'intervallo $[0.1, 0.01]$ (suff. piccolo).
- Criteri basati sulla generalizzazione:
L'algoritmo converge se le *performance di generalizzazione* sono adeguate.

Algoritmo di back-propagation:

- 1) **Inizializza** il vettore $\mathbf{w}(k)$ casualmente
- 2) **Presenta** gli esempi di training (per epoca o per pattern)
- 3) **Computa in avanti**, ottenendo l'output y_j e l'errore e_j per ogni j -esimo *neurone di output*
- 4) **Computa all'indietro** per aggiornare i pesi (differentemente a seconda che i neuroni siano *hidden* o *output*)
- 5) **Torna a (2)**, se la *condizione di terminazione* non è rispettata

Caratteristiche di una Rete Neurale:

- Rappresentazione dei dati (dominio delle feature)

Dipende dal problema, generalmente valori continui (reali), potenzialmente racchiusi in un intervallo limitato.

La **normalizzazione** è una trasformazione tale per cui le *feature* delle istanze del *training set* siano comprese nell'intervallo $[0, 1]$.

$$x_i = \frac{x_i - \min_i}{\max_i - \min_i}$$

- Topologia della rete (numero di livelli e di neuroni)

Dipende dal *task*. Nella pratica viene decisa tramite *trial and error*.

Per determinarla possono essere usati due tipi di algoritmi adattivi:

- **Pruning:** se le prestazioni sono soddisfacenti, si rimuovono neuroni e link finché continuano ad esserlo
- **Aggiunta di neuroni e link:** se le prestazioni non sono soddisfacenti, finché non lo diventano

- Parametri della rete (pesi, learning rate, numero di esempi del *training set*, ecc.)

- In generale, i **pesi iniziali** sono casuali; tipicamente con valori compresi negli intervalli $[-1, 1]$ e $[-0.5, 0.5]$.
- Il valore di η dipende dall'applicazione. Tipicamente nell'intervallo $[0.1, 0.9]$ oppure adattato durante l'allenamento.

- Apprendimento:

- **"Regola del pollice":** il numero di *istanze di training* dovrebbe essere da 4 a 10 volte il numero di pesi nella rete.
- **Altra regola:** $N > \frac{|W|}{(1-a)}$, dove $|W|$ è il numero di pesi e a è l'*accuratezza* attesa del test set.

- Approssimazione:

- **Funzioni booleane:** rappresentabili in un singolo *hidden layer*.
- **Funzioni continue limitate a tratti:** approssimabili con un errore arbitrariamente piccolo in un singolo *hidden layer*.
- **Funzioni continue:** approssimabili con un errore arbitrariamente piccolo in due *hidden layer*.

Teorema di esistenza dell'approssimazione universale:

Sia $\varphi(\cdot)$ una funzione continua non costante, limitata e monotona crescente.

Sia I_{m_0} un ipercubo unitario $[0, 1]^{m_0}$ m_0 -dimensionale.

Allora, data una funzione $f \in C(I_{m_0})$ e un $\varepsilon > 0$, esistono un intero m_1 e un insieme di costanti reali α_i , b_i e w_{ij} tali che:

$$F(x_1, \dots, x_{m_0}) = \sum_{i=1}^{m_1} \alpha_i \varphi\left(\sum_{j=1}^{m_0} w_{ij} x_j + b_i\right)$$

i. e. $|F(x_1, \dots, x_{m_0}) - f(x_1, \dots, x_{m_0})| < \varepsilon$

$F(x_1, \dots, x_{m_0})$ rappresenta l'*output* dell'*MLP* con:

- m_0 *nodi input*
- m_1 *nodi hidden*
- w_{ij} pesi sinaptici per *nodi hidden*
- b_i bias per *nodi hidden*
- α_i pesi sinaptici per *nodi output*

Osservazioni sul teorema: il teorema afferma che un singolo *neurone hidden* è sufficiente ad una *MLP* per computare un'approssimazione uniforme per un *training set* $\{x_1, \dots, x_{m_0}\}$.

Nel 1993, Barron ha stabilito la proprietà di approssimazione di *MLP*, valutando il **tasso di decrescita dell'errore** $O(1/m_1)$.

Le *FNN* possono essere applicate per risolvere problemi non linearmente separabili, per apprendere funzioni non lineari (regressione) e in particolare funzioni il cui input è una sequenza di misurazioni nel tempo (time series).

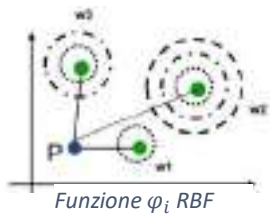
4 – RBF: Radial-Basis Function Networks

Modello di RN caratterizzato dalla funzione di attivazione lineare $\varphi(\mathbf{v}) = \mathbf{w}_1 \varphi_1(\|\mathbf{x} - \mathbf{t}_1\|) + \dots + \mathbf{w}_{m1} \varphi_{m1}(\|\mathbf{x} - \mathbf{t}_{m1}\|)$, che è la combinazione lineare di m funzioni di base radiale φ_i e dal fatto di avere **un solo hidden layer**.

Le Reti **RBF** sono utilizzate per task di regressione e classificazione complessa (non lineare).

Una rete RBF implementa l'approccio, differente dall'approssimazione stocastica delle reti MLP, chiamato **approccio ibrido**:

- 1) **Trasforma** un insieme di pattern non linearmente separabili in modo da renderli più probabilmente linearmente separabili.
- 2) **Classifica** l'insieme dei pattern trasformati.



Funzione φ_i RBF

$\|\mathbf{x} - \mathbf{t}_i\|$ è la distanza del vettore \mathbf{x} dal vettore \mathbf{t}_i

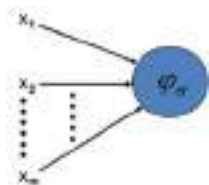
- I neuroni dell'hidden layer, usando una funzione RBF, applicano una trasformazione non lineare dall'input space all'hidden (o feature) space (il quale tipicamente è molto dimensionale).

In una funzione di base radiale (RBF), l'output è in funzione della distanza dal vettore prefissato \mathbf{t}_i .

- Il neurone dell'output layer combina linearmente l'output dei neuroni dell'hidden layer precedente e il peso di ciascuno di essi.

Teorema di Cover della separabilità dei pattern: il numero di neuroni dell'hidden layer e la probabilità che il problema sia linearmente separabile nel feature space, sono direttamente proporzionali.

Modello neurone hidden:



$\varphi_\sigma(\|\mathbf{x} - \mathbf{t}\|)$, dove \mathbf{t} e σ sono parametri:

\mathbf{t} è il parametro centro

σ è il parametro spread (ampiezza)

Se la funzione è gaussiana, la sensibilità del neurone può essere manipolata attraverso la variazione del parametro σ , il cui valore è inversamente proporzionale alla sensibilità.

Problema dell'interpolazione:

Dato un dataset di N diversi punti $\{\mathbf{x}_i \in \mathbb{R}^m, i = 1 \dots N\}$ e un vettore **desired output (target)** \mathbf{d} con N componenti reali $\{\mathbf{d}_i \in \mathbb{R}, i = 1 \dots N\}$, trovare una funzione $f: \mathbb{R}^m \rightarrow \mathbb{R}$ che soddisfi la **condizione di interpolazione**: $F(\mathbf{x}_i) = \mathbf{d}_i$

Se $F(\mathbf{x}) = \sum_{i=1}^N \mathbf{w}_i \varphi(\|\mathbf{x} - \mathbf{x}_i\|)$, allora:

$$\begin{bmatrix} \varphi_1(\|\mathbf{x}_1 - \mathbf{t}_1\|) & \dots & \varphi_{m1}(\|\mathbf{x}_1 - \mathbf{t}_{m1}\|) \\ \vdots & & \vdots \\ \varphi_1(\|\mathbf{x}_N - \mathbf{t}_1\|) & \dots & \varphi_{m1}(\|\mathbf{x}_N - \mathbf{t}_{m1}\|) \end{bmatrix} \begin{bmatrix} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_{m1} \end{bmatrix} = \begin{bmatrix} \mathbf{d}_1 \\ \vdots \\ \mathbf{d}_{m1} \end{bmatrix} \Rightarrow \Phi \cdot \mathbf{w}^T = \mathbf{d}$$

Teorema di Micchelli:

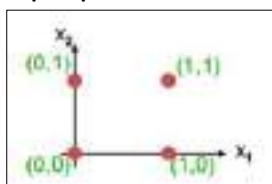
Sia $\{\mathbf{x}_i\}_{i=1}^N$ un insieme di punti $\in \mathbb{R}^m$ e distinti. Allora la **matrice Φ di interpolazione** tale per cui ciascuna componente ij sia $\varphi_{ij} = \varphi(\|\mathbf{x}_j - \mathbf{x}_i\|)$, se di dimensioni $N \times N$, è **non singolare**; ovvero il suo **determinante** è $\neq 0$ (il suo **rango non massimo**).

Tipi di φ per cui è valido il teorema di Micchelli:

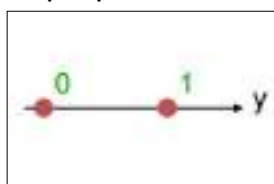
- **Multiquadratica:** $\varphi(r) = (r^2 + c^2)^{1/2}$
- **Multiquadratica inversa:** $\varphi(r) = \frac{1}{(r^2 + c^2)^{1/2}}$ $r = \|\mathbf{x} - \mathbf{t}\|$; $c > 0$
- **Gaussiana (normale):** $\varphi(r) = \exp\left(-\frac{r^2}{2\sigma^2}\right)$ $\sigma > 0$

Esempio: problema dello XOR

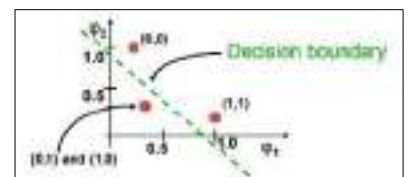
Input space:



Output space:



Mappati nello **spazio delle feature** $\langle \varphi_1, \varphi_2 \rangle$, C_1 e C_2 diventano **linearmente separabili** e il classificatore lineare che riceve gli input $\varphi_1(\mathbf{x})$ e $\varphi_2(\mathbf{x})$ può risolvere il problema.



Parametri di una Rete RBF necessari ad apprendere:

- 1) Il centro della funzione di attivazione RBF
- 2) L'ampiezza della funzione RBF, se Gaussiana
- 3) I pesi dall'hidden all'output layer

Stabilire i parametri di una rete RBF:

1) Metodo della matrice pseudo inversa:

- Centri: selezionati **casualmente** dal training set
- Ampiezze dei neuroni hidden: scelte attraverso **normalizzazione** $\sigma = \frac{\text{Maximum distance between any 2 centers}}{\sqrt{\text{number of centers}}} = \frac{d_{\max}}{\sqrt{m_1}}$

Funzione di attivazione del *neurone hidden* i : $\varphi_i(\|x - t_i\|^2) = \exp\left(-\frac{m_1}{d_{\max}^2} \cdot \|x - t_i\|^2\right)$

- Pesi: computabili attraverso l'**algoritmo della matrice pseudo-inversa**:

Per ogni istanza (x_i, d_i) , vorremmo che $\varphi(x_i) = d_i$.

Questo vincolo può essere imposto attraverso l'*equazione matriciale* $\Phi w = d$:

$$\begin{bmatrix} \varphi_1(\|x_1 - t_1\|) & \cdots & \varphi_{m_1}(\|x_1 - t_{m_1}\|) \\ \vdots & & \vdots \\ \varphi_1(\|x_N - t_1\|) & \cdots & \varphi_{m_1}(\|x_N - t_{m_1}\|) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_{m_1} \end{bmatrix} = [d_1 \quad \cdots \quad d_N]^T$$

Che può essere riscritta come $\Phi^T \Phi w = \Phi^T d$.

Matrice pseudo-inversa Φ^+ (o pseudo-inversa di Moore-Penrose):

generalizzazione della *matrice inversa* al caso in cui Φ non sia quadrata.

Definizione: $\Phi^+ \equiv (\Phi^T \Phi)^{-1} \Phi^T$.

Obiettivo: calcolare il *vettore pesi* che segue dall'*equazione matriciale* precedente $[w_1 \dots w_{m_1}]^T = \Phi^+ [d_1 \dots d_N]^T$

Problema: computare Φ^+ normalmente richiede di mantenere in memoria l'intera matrice Φ .

Idea: Dividere Φ in Q sottoinsiemi, ciascuno con un arbitrario numero a_i di righe ($1 \leq i \leq Q$).

Algoritmo:

- a. **Costruire** un insieme di nuove matrici A_i tali che ciascuna di esse sia composta da righe di zeri tali che: $\Phi = \sum_{i=1}^Q A_i$

- b. **Sostituire** questa definizione di Φ nella definizione di Φ^+ , ottenendo così:

$$w = (\Phi^T \Phi)^{-1} \Phi^T d = \left(\sum_{j,i=1}^Q A_j^T A_i \right)^{-1} \left(\sum_{j=1}^Q A_j^T \right) d$$

- c. Definire la matrice $H = \sum_{j,i=1}^Q A_j^T A_i$

Osservazioni:

- Tutti i prodotti $A_j^T \cdot A_i$ producono una matrice quadrata (ovvero $M \times M$)
- Per via della costruzione di A_i , tutti i prodotti $A_j^T \cdot A_i$ con $j \neq i$ producono una matrice *quadrata nulla*
- Sia \hat{A}_i un insieme di **matrici ridotte**, ovvero un insieme di matrici tali che ciascuna sia composta da una sola riga, quella non nulla, della matrice A_i da cui è generata. Otteniamo che $A_i^T A_i \equiv \hat{A}_i^T \hat{A}_i$.

- d. Computare $H^{-1} = \left(\sum_{j,i=1}^Q A_j^T A_i \right)^{-1}$:

1. Memorizzare \hat{A}_i (dimensioni di \hat{A}_i : $a_i \times M$)
2. Computare $\hat{A}_i^T \hat{A}_i$
3. Aggiungere il prodotto ottenuto al punto (2) a una lista v ;

Dopo Q iterazioni, $v = H$. Con quest'algoritmo viene mantenuta in memoria una sola \hat{A}_i per volta.

Inoltre, i prodotti possono essere computati in parallelo.

- e. Computare $w = H^{-1} \left(\sum_{j=1}^Q A_j^T \right) d$:

Sia Q l'insieme di **vettori ridotti** \hat{d}_i , tali per cui ciascun \hat{d}_i sia un insieme di componenti del vettore d

Analogamente al passo precedente vale $A_j^T d \equiv \hat{A}_j^T \hat{d}_j$ e quindi: $w = H^{-1} \left(\sum_{j=1}^Q \hat{A}_j^T \right) \hat{d}_j$

1. Memorizzare \hat{A}_j
2. Computare $\hat{A}_j^T \hat{d}_j$
3. Aggiungere il termine calcolato al punto (2) a una lista
4. Moltiplicare il termine ottenuto per H^{-1}

2) Hybrid Learning Process

- Centri: determinati con un algoritmo di clustering
 - a. Inizializzazione: inizializza un vettore \mathbf{t} le cui m_1 componenti (inizialmente casuali) sono cluster
 - b. Campionamento: rappresenta \mathbf{x} dallo spazio dell'input
 - c. Similarity matching: trova l'indice del cluster più vicino $k(\mathbf{x}) = \arg \min_k \|\mathbf{x}(n) - \mathbf{t}_k(n)\|$
 - d. Aggiornamento: avvicina $\mathbf{t}_{k(\mathbf{x})}$ ad \mathbf{x} di una misura pari al *learning rate* η ponderato per la distanza di \mathbf{x} da \mathbf{t}

$$\mathbf{t}_k(n+1) = \begin{cases} \mathbf{t}_k(n) + \eta[\mathbf{x}(n) - \mathbf{t}_k(n)] & \text{se } k = k(\mathbf{x}) \\ \mathbf{t}_k(n) & \text{altrimenti} \end{cases}$$

- e. Continuazione: finché vi sono cambiamenti del centro consistenti, incrementa n di 1 e ricomincia dal punto (2).
- Ampiezze dei neuroni hidden: scelte attraverso **normalizzazione**
 - Pesi: computati attraverso l'**algoritmo LMS**

3) Metodo della discesa del gradiente per trovare i centri, lo spread e i pesi, minimizzando l'errore quadratico istantaneo

- Centri: $\Delta \mathbf{t}_j = -\eta_{t_j} \frac{\partial E}{\partial \mathbf{t}_j}$
- Ampiezze dei neuroni hidden: $\Delta \sigma_j = -\eta_{\sigma_j} \frac{\partial E}{\partial \sigma_j}$
- Pesi: $\Delta \mathbf{w}_{ij} = -\eta_{w_{ij}} \frac{\partial E}{\partial \mathbf{w}_{ij}}$

Comparazione tra **RBF** e **FFNN**:

- Entrambe le architetture sono **stratificate in modo non lineare**
- Entrambe sono considerabili **approssimatori universali**

Architettura:

- Reti **RBF** hanno un **singolo hidden layer**
- Reti **FFNN** possono avere **molteplici hidden layer**

Modello dei neuroni:

- Nelle Reti **RBF** i **modelli** dei neuroni *hidden* sono **diversi** dal modello del neurone in *output*
- Nelle Reti **FFNN** neuroni *hidden* e neurone di *output* tipicamente condividono un **modello comune**
- Nelle Reti **RBF** gli **hidden layer** sono **non lineari** e l'**output layer** è **lineare**
- Nelle Reti **FFNN** tipicamente sia gli **hidden layer** che l'**output layer** sono **non lineari**

Funzione di attivazione dei neuroni hidden:

- Nelle Reti **RBF** computa la **distanza euclidea** tra i vettori in input e il centro di quella unità.
- Nelle Reti **FFNN** computa il **prodotto interno** dei vettori in input e il vettore peso di quel neurone

Approssimazione:

- Nelle Reti **RBF** la funzione gaussiana costruisce un'**approssimazione locale** dell'associazione non lineare I/O
- Nelle Reti **FFNN** costruisce un'**approssimazione globale** dell'associazione non lineare I/O

5 – ELM: Extreme Learning Machine

Rete SLF: *Single hidden Layer Feedforward Neural Network*

ELM: algoritmo di apprendimento per *Reti SLF* con **prestazioni computazionali** fino a centinaia di volte superiori dei tradizionali algoritmi di apprendimento basati sulla *backpropagation* e con **prestazioni di generalizzazione** migliori.

Definizione: una funzione si dice *infinitamente differenziabile* in un punto se la sua n -esima derivata esiste in quel punto $\forall n \in \mathbb{N}$.

Teoremi principali di ELM:

Ipotesi:

1. Sia data una *Rete SLF* con N neuroni *hidden* e sia $g: \mathbb{R} \rightarrow \mathbb{R}$ la loro *funzione di attivazione infinitamente differenziabile*.
2. Siano N un numero arbitrario di istanze distinte $(\mathbf{x}_i, \mathbf{t}_i)$, dove $\mathbf{x}_i \in \mathbb{R}^n$ e $\mathbf{t}_i \in \mathbb{R}^m$, e \mathbf{T} il vettore di componenti \mathbf{t}_i .
3. Siano $\boldsymbol{\beta}_i$ il **vettore dei pesi** che associa l' i -esimo neurone *hidden* ai neuroni *output* e $\boldsymbol{\beta}$ il vettore di componenti $\boldsymbol{\beta}_i$.
4. Sia ε un qualsiasi numero arbitrariamente piccolo positivo.

1) Allora, per ogni $\mathbf{w}_i \in \mathbb{R}^n$ tra neuroni *input* ed *hidden* e $b_i \in \mathbb{R}$ scelti casualmente in funzione di una certa distribuzione di probabilità continua, è garantito che esista la matrice \mathbf{H} *output* dell'*hidden layer* invertibile tale che valga $\|\mathbf{H}\boldsymbol{\beta} - \mathbf{T}\| = 0$.

2) Allora, esiste $\tilde{N} \leq N$ tale che $|\mathbf{H}_{N \times \tilde{N}} \cdot \boldsymbol{\beta}_{\tilde{N} \times m} - \mathbf{T}_{N \times m}| < \varepsilon$, dove \tilde{N} è il numero di neuroni *hidden*

Quindi: se le *funzioni di attivazione* dell'*hidden layer* di una sono *infinitamente differenziabili*, allora:

- 1) I **vettori pesi input** e **hidden layer** possono essere inizializzati casualmente
- 2) La *Rete SLF* può essere considerata come un sistema lineare e il **vettore pesi output** $\boldsymbol{\beta}$ può essere determinato analiticamente attraverso l'**operazione inversa generalizzata** \mathbf{H}^+ delle matrici di *output* dell'*hidden layer* \mathbf{H} .
Questo è molto utile in quanto non sempre \mathbf{H} è quadrata.

Sia $\varphi_i(v) = \sum_{h=1}^{\tilde{N}} \boldsymbol{\beta}_h \cdot g(\mathbf{x}_j)$ la funzione di attivazione del i -esimo neurone *output* in risposta al j -esimo *input*.

Sebbene sia possibile il calcolo di $\boldsymbol{\beta}$ in funzione di \mathbf{H}^{-1} , non è quello a cui siamo interessati poiché \mathbf{H} non è invertibile per $\tilde{N} < N$.

$$\begin{bmatrix} g(\bar{\mathbf{x}}_1 \cdot \bar{\mathbf{w}}_1) & \cdots & g(\bar{\mathbf{x}}_1 \cdot \bar{\mathbf{w}}_{\tilde{N}}) \\ \vdots & & \vdots \\ g(\bar{\mathbf{x}}_{\tilde{N}} \cdot \bar{\mathbf{w}}_1) & \cdots & g(\bar{\mathbf{x}}_{\tilde{N}} \cdot \bar{\mathbf{w}}_{\tilde{N}}) \end{bmatrix} \begin{bmatrix} \beta_{11} \cdots \beta_{1m} \\ \vdots \\ \beta_{\tilde{N}1} \cdots \beta_{\tilde{N}m} \end{bmatrix} = \begin{bmatrix} y_{11} \cdots y_{1m} \\ \vdots \\ y_{\tilde{N}1} \cdots y_{\tilde{N}m} \end{bmatrix}$$

$$\Rightarrow \mathbf{H} \cdot \boldsymbol{\beta} = \mathbf{y} \Rightarrow \boldsymbol{\beta} = \mathbf{H}^{-1} \cdot \mathbf{y}$$

Bensì vogliamo risolvere la seguente equazione matriciale: $\mathbf{H} \cdot \boldsymbol{\beta} \approx \mathbf{T} \Rightarrow \boldsymbol{\beta} = \mathbf{H}^+ \cdot \mathbf{T}$, che è invece sempre risolvibile (ma $\mathbf{T} \approx \mathbf{y}$).
Ovvero effettuare **regressione** e non **interpolazione**.

Algoritmo ELM:

1. Assegna valori casuali a \mathbf{w}_i e b_i , $i = 1, \dots, \tilde{N}$ (l'*upper bound* del numero di nodi richiesti è N)
2. Computa \mathbf{H}
3. Computa la minima *norma least-square* $\boldsymbol{\beta} = \mathbf{H}^+ \mathbf{T}$

L'algoritmo converge per ogni funzione di attivazione g *infinitamente differenziabile*.

Minore è la norma del vettore pesi, maggiori sono le performance di generalizzazione.

Questo metodo è affetto da severe limitazioni quando la matrice è quasi *singolare*.

Matrice singolare: matrice *quadrata* ma *non invertibile*, poiché il suo *determinante* è $= 0$ (ovvero il suo *rango* non è pieno)

Approcci:

- 1) Aggiunta di un termine di regolarizzazione, così che il costo funzionale sia esprimibile come:

$$E = E_D + \lambda E_W = \frac{1}{2} \sum_{j=1}^N \sum_{k=1}^Q \left(\left(\mathbf{t}_j^k - \sum_{i=1}^M \mathbf{w}_{ki} \phi(c_i \cdot \mathbf{x}_j + b_i) \right)^2 + \frac{\lambda}{2} \sum_{i=1}^M |\mathbf{w}_{ki}|^2 \right)$$

La cui soluzione è: $\mathbf{w}^* = (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}^T$

2) Calcolo della Singular Value Decomposition di \mathbf{H} :

Metodo computazionalmente più efficace. Inoltre, poiché non è necessario il *termine di regolarizzazione*, è possibile realizzare una *minimizzazione dell'errore dipendente dai dati* più efficace.

Definizioni:

- **Matrice unitaria:** matrice \mathbf{U} che soddisfa la condizione $\mathbf{U}^\dagger \mathbf{U} = \mathbf{U} \mathbf{U}^\dagger = \mathbf{I}$
- **Matrice diagonale:** matrice in cui solamente i valori della diagonale principale possono essere diversi da 0.
- **Matrice trasposta coniugata:** matrice \mathbf{V} che soddisfa la condizione $\mathbf{V}^\dagger = (\mathbf{V}^T)^* = (\mathbf{V}^*)^T$

Singular Value Decomposition (SVD): fattorizzazione di una matrice basata sull'uso di **autovettori**.

Data una matrice \mathbf{M} di dimensioni $m \times n$, si tratta di una scrittura del tipo: $\mathbf{M} = \mathbf{U} \mathbf{\Sigma} (\mathbf{V}^T)^* = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\dagger$, dove:

- \mathbf{U} è una **matrice unitaria** di dimensioni $m \times m$
- $\mathbf{\Sigma}$ è una **matrice diagonale** di dimensioni $m \times n$ le cui componenti vengono indicate come σ_{ij}
- \mathbf{V}^\dagger è la matrice **trasposta coniugata** di una **matrice unitaria** \mathbf{V} di dimensioni $n \times n$

Le componenti sulla diagonale di $\mathbf{\Sigma}$ sono detti **valori singolari** di \mathbf{M} .

Le colonne di \mathbf{U} sono dette **vettori singolari sinistri** di \mathbf{M} (o **autovettori** di $\mathbf{M} \mathbf{M}^*$).

Le colonne di \mathbf{V} sono dette **vettori singolari destri** di \mathbf{M} (o **autovettori** di $\mathbf{M}^* \mathbf{M}$).

Usando questa strategia, la **matrice pseudoinversa** \mathbf{H}^+ è definita come: $\mathbf{H}^+ = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^T$, ovvero evitando l'inversione *matriciale*.

- $\mathbf{\Sigma}^+$ è una **matrice diagonale** le cui componenti *non nulle* assumono il valore $\frac{1}{\sigma_{ij}}$

L'inversione *numerica* di valori molto bassi σ_{ij} in $\mathbf{\Sigma}$ potrebbe rendere instabile l'algoritmo.

6 – CNN: Convolutional Neural Networks

Deep Neural Network (DNN): reti caratterizzate da un'ampia profondità e da una struttura gerarchica di layer.

- La **profondità** (numero di *layer*) è solo un fattore di complessità. **Neuroni, connessioni** e conseguentemente **numero di pesi** caratterizzano anche la complessità di un **DNN**.
- L'**organizzazione gerarchica** consente di condividere e riutilizzare le informazioni. Lungo la gerarchia è possibile selezionare una caratteristica specifica e scartare i dettagli non necessari (al fine di massimizzare l'*invarianza*).

Cenni storici

Per motivi di **efficienza**, gli studi inerenti le reti neurali per molto tempo si sono concentrati su reti con un singolo *hidden layer*: funzioni computabili con complessità polinomiale su n layer possono richiedere complessità esponenziale con un differente numero di layer. Questo perché all'aumentare del numero di pesi presenti nella rete, aumentano anche le risorse computazionali necessarie al suo addestramento.

Principali modelli DNN:

Addestramento supervisionato, modelli discriminativi:

- **CNN - Convolutional Neural Networks**
- **FC DNN - Fully Connected DNN (MLP with at least two hidden layers)**
- **HTM - Hierarchical Temporal Memory**

Addestramento supervisionato, modelli generativi (utili per il pre-training):

- **Stacked (de-noising) Auto-Encoders**
- **RBM - Restricted Boltzmann Machine**
- **DBN - Deep Belief Networks**

Modelli ricorrenti,

- **RNN - Recurrent Neural Network**
- **LSTM - Long Short-Term Memory**

Motivi del successo delle DNN:

Le **CNN** sono state introdotte nel 1998 e hanno ottenuto un modesto successo per problemi di piccole dimensioni.

Nel 2012, con l'arrivo della rete convoluzionale **AlexNet**, si è verificato un cambiamento radicale.

AlexNet non introdusse innovazioni importanti rispetto alla **CNN**, ma alcuni fattori ne hanno decretato il successo:

- o La presenza di **BigData** (*dataset* con immense quantità di istanze etichettate)
- o L'introduzione delle **unità di calcolo grafico (GPU)**: architetture che a differenza delle *unità di calcolo centrali (CPU)* favoriscono l'elaborazione di calcoli paralleli e che hanno drasticamente ridotto i tempi necessari all'apprendimento.
- o La sostituzione della funzione di attivazione *sigmoide* in favore della **Relu**.

Convoluzione: operazione matematica attraverso la quale vengono applicati *filtri* ad una matrice.

L'applicazione viene effettuata, iterativamente e per scorrimento, a diverse porzioni dell'input attraverso il prodotto scalare.

Filtro: *maschera tensoriale di pesi*.

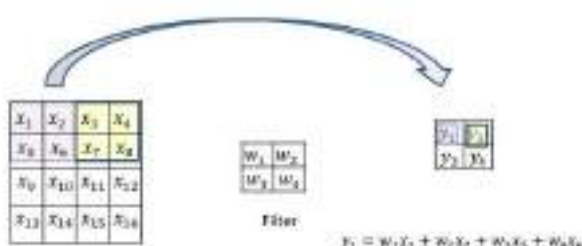
Architettura di una CNN: gerarchia di layer

- **Convolutional layers:** sequenza di *layer* che esegue la scansione dell'input applicando *filtri di convoluzione*.

L'*output* prende il nome di *feature map*. Le dimensioni della *feature map* dipendono da:

- o **Dimensioni del filtro:** la profondità della *feature map* è uguale alla profondità del filtro.
- o **Stride:** numero di unità di cui si muove per scorrimento la finestra del filtro.
- o **Zero-padding:** presenza di vettori colonna o riga nulli ai bordi della matrice di *input*.

Esempio:



- **Pooling layers:** sequenza di layer che esegue *downsampling* (approssimazione di una sequenza in *input*).

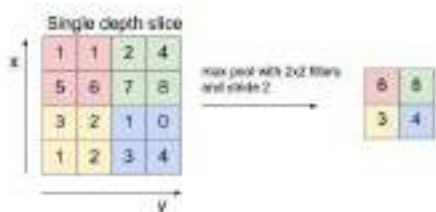
Il layer di pooling esegue aggregazione di informazioni, generando *feature map* di dimensioni minori dell'*input*.

Obiettivo: trasformare l'*input* mantenendo *invarianza* (ovvero mantenendo solo informazioni significative per la discriminazione).

Esempi di pooling: pooling massimo e pooling medio; prendono rispettivamente valore massimo e valore medio della sequenza. Gli operatori che lo effettuano, *avg* e *max*, prendono il nome di *operatori aggregati*.

Anche in questo caso l'*output* prende il nome di *feature map*; le sue dimensioni dipendono dalle dimensioni del filtro.

Questo tipo di aggregazione non necessita di parametri/pesi per apprendere.



Nell'algoritmo di *backpropagation* i gradienti ai vari livelli vengono moltiplicati tramite la *chain rule*.

Scomparsa del gradiente: gli n fattori $\in (0,1) \rightarrow$ il prodotto decresce *esponenzialmente* rispetto a n (profondità della rete).

Esplosione del gradiente: gli n fattori assumono valori elevati \rightarrow il prodotto decresce *esponenzialmente* rispetto a n .

- **Fully connected layers:** sequenza di layer che compone una tradizionale Rete MLP.

Funzione di attivazione:

- o **Relu (Rectified Linear):** funzione di attivazione che risolve i problemi di scomparsa ed esplosione del gradiente. La derivata è 0 per valori negativi, 1 per quelli positivi. Questo porta ad *attivazioni sparse* e *robustezza* della rete.

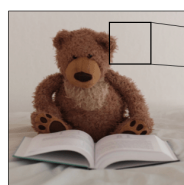
Scelte comuni nelle architetture CNN moderne:

- o **Soft-max layer:** ultimo livello della rete, *completamente connesso* al precedente. Ciascun neurone i computa una funzione logistica generalizzata: prende come *input* il vettore di punteggi $x \in \mathbb{R}^n$ generati dal livello precedente e genera un vettore di probabilità $p \in \mathbb{R}^n$, dove n è il numero di neuroni.

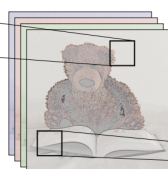
$$p_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

- o **Cross-entropy loss function:** funzione che quantifica la differenza tra due distribuzioni p e q .

$$H(p, q) = - \sum_v p(v) \cdot \log(q(v))$$



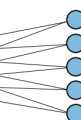
Input image



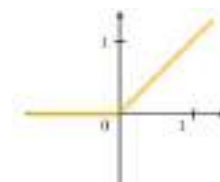
Convolutions



Pooling



Fully Connected



Funzione di attivazione Relu
 $\varphi = \max(0, z)$

Caratteristiche di CNN:

Local Processing:

- Ciascun neurone effettua un processamento locale, ovvero differente da quello degli altri neuroni allo stesso livello.

Shared weights:

- Gli strati di una *feature map* condividono lo stesso *vettore pesi*.

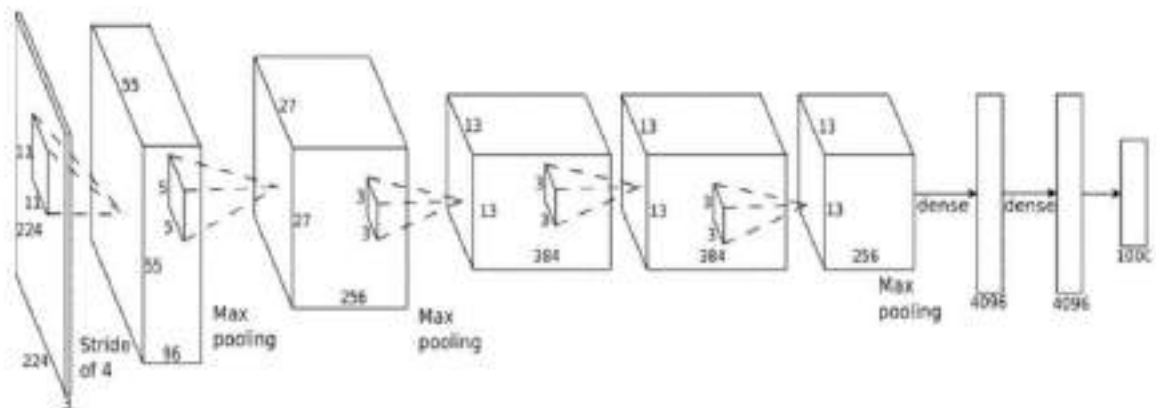
Insieme hanno contribuito a una forte riduzione del numero di pesi necessari alla rete per fare apprendimento.

Allenamento di una CNN: molto lento, a causa della necessità di dataset molto grandi.

Fortunatamente, una volta terminato l'apprendimento della rete, apprendere nuovi *pattern* è generalmente un'operazione più veloce. Il problema rimane la possibilità di usufruire di un dataset di dimensioni opportune.

Transfer learning: memorizzare le conoscenze acquisite durante la risoluzione di un problema e applicarle a un problema diverso.

- **Fine-tuning:** apprendimento di un modello a partire da una *rete precedentemente allenata* su un problema simile:
 - 1) **Sostituire** l'*output layer* con un nuovo *layer* (tipicamente aggiornando il numero delle classi).
 - 2) **Conservare** i pesi della *rete* già allenata, tranne quelli agli ultimi due *layer*.
 - 3) **Effettuare** ulteriori iterazioni di apprendimento per ottimizzare i pesi in rispetto alle peculiarità del nuovo dataset
- **Features re-using:** a partire da una *rete precedentemente allenata* (su cui non è stato effettuato *fine-tuning*):
 - 1) **Estrarre** da un *layer intermedio* le *feature* generate dalla *rete*
 - 2) **Usare** tali *feature* per allenare un classificatore esterno (ad esempio **SVM**) a classificare nuovi pattern.



7 – SOM: Self-Organizing Maps

Modello di Rete Neurale caratterizzato un tipo di apprendimento **non supervisionato**: l'**apprendimento competitivo**, il quale prevede che i neuroni di *output* competano tra loro per essere attivati.

Il neurone di *output* che vince la competizione prende il nome di *neurone vincitore*.

Nelle *SOM* i neuroni sono posti in un *reticolo* che definisce un **ordine topologico**, tipicamente **mono/bi-dimensionale**.

Tali neuroni si **regolano selettivamente** alla ricezione di **stimoli** (ovvero di *pattern input*), la quale avviene in modo continuo fino al termine dell'apprendimento. Tali *stimoli* modificano la posizione dei neuroni, i quali catturano le regolarità delle *feature* e producono un sistema di coordinate significativo rispetto alla loro varianza che prende il nome di **feature map**.

"The brain is organized in such a way that different sensory inputs are represented by topologically ordered maps. Sensory inputs such as tactile, visual and acoustic inputs are mapped onto different cerebral cortex areas."

Ciascun *input pattern* presentato alla rete produce l'attivazione di una regione localizzata della *feature map*.

La posizione e la natura di dell'attivazione sono in funzione delle *feature* dell'*input pattern*.

Affinché l'addestramento della rete si sviluppi correttamente, essa deve essere stimolata attraverso input sufficientemente diversi tra loro. La presenza di ridondanza è fondamentale, infatti gli stimoli in fase di apprendimento vengono campionati con una certa distribuzione di probabilità a partire dall'*input space*.

A ciascuna *feature* dell'*input pattern* x è connesso ciascun neurone j della rete, al quale è associato a un vettore pesi w_j .

Processi di una SOM:

- 1) **Competizione:** per ogni *input pattern* x , i neuroni nella rete computano i rispettivi valori di una *funzione discriminante* $i(x)$.

Tale *funzione discriminante*, $i(x) = \arg \min_j \|x - w_j\|$, designa il **neurone vincitore** i (o **Best Matching Unit, BMU**).

Minimizzare la norma 2 (distanza euclidea) e massimizzare il prodotto scalare sono criteri che coincidono.

Selezione del neurone vincitore:

1. La rete viene stimolata da x
2. Per ogni j , la rete calcola $w_j^T x$
3. Poi seleziona il neurone associato al prodotto maggiore (distanza euclidea minore).
Tale neurone sarà il centro di un **vicinato topologico** di *neuroni eccitati*.

- 2) **Cooperazione:** il *neurone vincitore* determina la posizione spaziale del **vicinato topologico** di *neuroni eccitati*.

Lateral interaction:

"In the human brain a neuron that is firing tends to excite the neurons in its immediate neighborhood more than those farther away from it, which is intuitively satisfying."

Sulla base di questo principio è stata modellata la **selezione del vicinato topologico**, centrato nel *neurone vincitore* i circondato da un insieme di *neuroni eccitati* j .

Una buona scelta (non l'unica possibile) di *funzione di selezione del vicinato topologico* h_{ji} è la funzione gaussiana:

$$h_{j,i(x)} = \exp\left(-\frac{d_{j,i}^2}{2\sigma^2}\right)$$

$d_{j,i}$ è la *distanza laterale* tra i e il neurone vincitore j , ovvero la distanza nel *feature map* (**non** nello spazio dell'*input*).

σ è il parametro che regola la **larghezza effettiva** del *vicinato topologico*.

- 3) **Adattamento sinaptico:** meccanismo che, aggiustando i *pesi sinaptici* w_j dei *neuroni eccitati*, incrementa i valori della *funzione discriminante* in relazione all'*input pattern* x . Tali aggiustamenti enfatizzano la risposta del *neurone vincitore* e del *vicinato topologico* a futuri stimoli simili ad x .

$$\begin{aligned}\hat{w}(n+1) &= \hat{w}_j(n) + \eta(n) h_{j,i(x)}(n) [x(n) - w_j(n)] \\ \Rightarrow \Delta w &= w(n+1) - w(n) = \eta(n) h_{j,i(x)}(n) [x(n) - w_j(n)] = \eta \cdot h_{j,i(x)} \cdot (x - w_j)\end{aligned}$$

Fasi del processo adattivo:

- 1) **Fase di ordinamento**, o di **auto-organizzazione** (~ 1000 iterazioni)

Fase in cui avviene l'*ordinamento topologico* dei vettori peso e in cui i parametri η e σ_0 assumono i seguenti valori:

$$\eta = 0.1 \text{ (lower bound di decrescita: 0.01)} \quad \sigma_0 = \text{raggio del reticolo}$$

- 2) **Fase di convergenza** ($\sim [500 * \# \text{neuroni}]$ iterazioni)

Fase in cui aumenta la precisione di corrispondenza tra *pattern input* e neuroni e in cui i parametri η e σ si stabilizzano intorno a $\eta = 0.1$ e σ = valore tale che $h_{j,i(x)}$ contenga solo i neuroni immediatamente più vicini al **BMU**.

Algoritmo di apprendimento:

- 1) Inizializza il **vettore pesi**, desiderabilmente con valori bassi e diversi per ciascun neurone.
- 2) **Campiona** uno *stimolo* a partire dall'*input space* seguendo una certa distribuzione di probabilità*.
 1. **Similarity matching**: ricerca il *neurone vincitore (BMU)*.
 2. **Aggiorna** i vettori peso dei *neuroni eccitati*.
- 3) **Torna a (2)**, se ci sono stati cambiamenti significativi nell'aggiornamento dei pesi.

L'algoritmo di apprendimento prevede un **decadimento esponenziale iterativo** in n di σ e η in funzione del *parametro temporale* τ

$$\eta(n) = \eta_o \exp\left(-\frac{n}{\tau}\right) \quad \sigma(n) = \sigma_o \exp\left(-\frac{n}{\tau}\right)$$

* La distribuzione di probabilità del campionamento influenza la *feature map*, all'interno della quale gli stimoli più frequenti saranno *meglio* ("con una maggiore risoluzione") rappresentati.

8 – Hopfield & Boltzman Networks

Energy based model: modello che si avvale di una *funzione*, detta di *energia globale*; tale funzione, dati due neuroni i, j e il peso w_{ij} della loro connessione, associa un valore (*energia*) alla loro specifica configurazione di stati.

Le Reti di Hopfield sono il più semplice tipo di *energy based model* e sono composte da **unità (neuroni) binarie di soglia**, collegate tra loro da **connessioni ricorrenti (bidirezionali)**. Tali reti sono caratterizzate dalla seguente **funzione di energia globale**:

$$E = - \sum_i s_i b_i - \sum_{i < j} s_i s_j w_{ij}$$

dove b_i è *bias* di ciascuna unità i , s_i è il suo *stato* e w_{ij} è il *peso del collegamento* tra le unità i e j .

Ciascuna unità i può **quantificare localmente** la propria influenza sull'*energia globale*, ovvero il proprio **energy gap**:

$$v_i(n) = \Delta E_i = E(s_i = -1) - E(s_i = 1) = b_i + \sum_j s_j w_{ij}$$

Inoltre, la *funzione di attivazione* φ è: $\varphi(v_j(n)) = \begin{cases} 1 & \text{if } v_j(n) > 0 \\ -1 & \text{if } v_j(n) < 0 \end{cases}$

Possibili outcome di reti ricorrenti di unità non lineari:

- Raggiungimento di uno stato stabile
- Oscillazione
- Perseguimento di traiettorie caotiche (non prevedibili)

Analisi del comportamento di una rete ricorrente con connessioni simmetriche e opportune funzioni di attivazione:

- Emerge una *funzione di energia globale*.
- **La rete si stabilizza nei punti di minimo** della *funzione di energia globale*.

Il tipico task che le Reti di Hopfield sono in grado di eseguire è il **pattern completion**: data una porzione sufficiente di informazione memorizzata, recuperarne la completezza. Lo realizzano attraverso un'operazione che prende il nome di **information withdrawal**.

Algoritmo non deterministico di information withdrawal (attraverso la minimizzazione della funzione di energia):

- 1) Calcola l'*energia globale* della configurazione di stati.
- 2) Applica ad un'unità casuale la *funzione di attivazione*.
- 3) Torna a (1), se la rete non si è *stabilizzata* (ovvero se la sua *configurazione* non corrisponde a un punto di minimo).

Teorema di convergenza: l'algoritmo **converge sempre**.

Dimostrazione in assenza di bias:

Sia N il numero di unità della rete. Allora esistono 2^N possibili configurazioni della rete.

A ciascuna configurazione, attraverso la *funzione di energia senza bias* $E = -\frac{1}{2} \sum_i \sum_j s_i s_j w_{ij}$, viene assegnato un valore.

Osservazione: $(s_i, s_j \text{ concordi} \wedge w_{ij} > 0) \vee (s_i, s_j \text{ discordi} \wedge w_{ij} < 0) \Rightarrow E > E'$, ovvero l'*energia globale* *decresce*.

Siano E ed E' due configurazioni successive nell'apprendimento: $E - E' = -\sum_{j \neq k} s_j (s_k - s'_k) w_{kj}$, k unità aggiornata.

- $s_k = 1 \wedge s'_k = -1 \Rightarrow s_k - s'_k > 0 \wedge \sum_j w_{kj} s_j < 0$
- $s_k = -1 \wedge s'_k = 1 \Rightarrow s_k - s'_k < 0 \wedge \sum_j w_{kj} s_j > 0$

Essendo il numero di possibili configurazioni della rete finito, è dimostrato che l'algoritmo converge.

Osservazioni:

1. La *funzione di energia globale* può avere più di un punto di minimo (o **memoria stabile**), ma non è garantito che l'algoritmo converga nel punto di **minimo assoluto** (o **memoria fondamentale**).
2. La regola di aggiornamento non è *error-driven*.

Principio di Hebb per l'apprendimento:

Rafforzare le connessioni tra unità aventi stesso stato e indebolire le connessioni tra unità con stato opposto.

Controparte biologica: il cervello

- **Rafforza** le sinapsi tra unità spesso attive allo stesso tempo.
- **Indebolisce** quelle tra i neuroni che non sono attivi simultaneamente.

Algoritmo di apprendimento:

Memorizzazione: date M memorie fondamentali, **calcola** ogni collegamento w_{ji} tale che $j \neq i$: $w_{ji} = \frac{1}{M} \sum_{k=1}^M f_k(i) f_k(j)$

Qualità: *pattern completion*, generalizzano input simili, permettono l'estrazione di prototipi, fault tolerant, Hebb rule, context effect.

Difetti: minimi spuri, sinapsi simmetriche e stati stabili (non corrisposti biologicamente - *stati transitori*), capacità limitata.

Capacità limitata delle reti di Hopfield: N unità $\rightarrow 0.15N$ memorie \rightarrow uso non efficiente della memoria dovuto al comportamento della funzione di energia globale, la quale **fonde due minimi** se la loro distanza è inferiore a una certa soglia.

Modifiche che migliorano la capacità delle reti di Hopfield:

- **Aggiunta di unità hidden alla Rete di Hopfield standard:**

Approccio teorizzato da Hopfield, Feinstein e Palmer, che però non ne hanno approfondito l'analisi e ispirato agli studi di Crick e Mitchison sull'importanza della fase REM nel sonno.

Realizzare una fase di "*unlearning*" a seguito di apprendimento al fine di superare eventuali *minimi locali*.

- Apprendimento: determinazione di nuovi minimi della funzione di energia globale.
- Unlearning: modifica dell'apprendimento effettuato.

Le **unità hidden** rappresentano l'interpretazione percettiva che la rete fa dell'input.

I **pesi** rappresentano i vincoli sulle interpretazioni.

La **funzione d'energia globale** misura la bontà dell'interpretazione.

- **Pseudo-likelihood technique:**

Approccio, scoperto dalla fisica Elizabeth Gardner, che consiste nel sostituire la regola di aggiornamento standard con l'applicazione della regola di convergenza del perceptrone su ciascuna unità della rete.

L'apprendimento con quest'approccio:

1. **Evita minimi locali** (ovvero *soluzioni sub-ottime*) con l'introduzione di *non determinismo*.

Simulated annealing (Kirkpatrick): euristica che introduce nell'apprendimento un *fattore di rumore (noise)*.

La sua implementazione richiede che le **unità binarie** siano **stocastiche** e non più **di soglia** (che invece sono *deterministiche*).

Stocasticamente ad ogni iterazione decide con probabilità p se transire dalla *configurazione di stati* attuale.

Probabilità di transizione: $p_j = \frac{1}{1 + e^{-\Delta E_j/T}}$

Osservazioni:

- p è una *funzione logistica normale*
- T (**temperatura**) è un parametro di *regolarizzazione*

2. **Massimizza la verosimiglianza (likelihood) congiunta** delle istanze v del *training set*, a ciascuna delle quali il modello assegna la verosimiglianza p_v :

$$\max_w \prod_v p_v \equiv \min_w \sum_v \log p_v$$

$$\frac{\partial \log(p_v)}{\partial w_{ij}} = \langle s_i \cdot s_j \rangle_v - \langle s_i \cdot s_j \rangle_{model} \Rightarrow \Delta w_{ij} \propto \langle s_i \cdot s_j \rangle_{data} - \langle s_i \cdot s_j \rangle_{model}$$

$\langle s_i \cdot s_j \rangle_{data}$ è il **termine di apprendimento**, ovvero il **valore atteso** di $s_i \cdot s_j$ dato dal vettore **data** (*media del training set*).

$\langle s_i \cdot s_j \rangle_{model}$ è il **termine di unlearning**, ovvero **valore atteso** di $s_i \cdot s_j$ quando la rete non viene sottoposta ad *input*.

Osservazione: tale l'algoritmo di allenamento eredita la stessa inefficienza dell'algoritmo del perceptrone.

Restricted Boltzmann Machines (RBM): semplificazione delle tradizionali *Reti di Boltzmann*.

Sono caratterizzate da **un solo hidden layer**, e dal fatto che le **unità** presenti sullo **stesso livello non sono collegate tra loro**.

Algoritmo di apprendimento di una RBM:

- 1) Sottopone l'*RBM* all'*input*
- 2) Attende che la *rete* raggiunga l'equilibrio termico, ovvero che la configurazione di stati non transisca più
- 3) Applica la regola di aggiornamento $\Delta w_{ij} = \epsilon(\langle s_i \cdot s_j \rangle^0 - \langle s_i \cdot s_j \rangle^\infty)$

Tale algoritmo è di per sé più efficace dell'algoritmo di apprendimento delle *Boltzman Machine*, poiché l'indipendenza dei neuroni sullo stesso livello permette di aggiornarli in parallelo. Esiste tuttavia una versione ulteriormente migliorata di tale algoritmo:

Contrastive Divergence Algorithm:

- 1) Computa $\langle s_i \cdot s_j \rangle^0$ e $\langle s_i \cdot s_j \rangle^1$, rispettivamente **data** e **reconstructor**.
- 2) Applica la regola di aggiornamento $\Delta w_{ij} = \epsilon(\langle s_i \cdot s_j \rangle^0 - \langle s_i \cdot s_j \rangle^1)$, che **non** segue il *gradiente della verosimiglianza* ma che è comunque efficace in quanto sfrutta la differenza di energia tra i due addendi (temperatura bassa nel primo, alta nel secondo).

Osservazione: i pesi catturano *feature globali* del *training set*, quindi la ricostruzione a partire da classi non familiari finisce comunque con l'obbedire alle regolarità delle feature catturate.