



卷一

目錄

卷二

卷三

“Si sta come d’autunno

sugli alberi

le foglie.”

*Se non puoi passare attraverso una montagna, giraci intorno;
se non puoi girarci intorno, passaci sopra;
se non puoi passarci sopra, siediti un attimo
e chiediti se raggiungere l'altro lato sia davvero così importante.*

Se lo è comincia a scavare una galleria.

Indice

| | |
|---|-----------|
| Introduzione | 5 |
| 1 Presentazione delle reti neurali artificiali | 6 |
| 1.1 Cenni storici ed origini | 6 |
| 1.2 Analogie con il sistema nervoso ed applicazioni delle ANN . . | 9 |
| 1.3 Il neurone biologico | 11 |
| 1.4 Modello del neurone artificiale | 13 |
| 1.5 Funzioni di attivazione | 17 |
| 1.5.1 Funzione a soglia | 17 |
| 1.5.2 Funzione lineare | 19 |
| 1.5.3 Funzione lineare a tratti | 19 |
| 1.5.4 Funzione sigmoide | 20 |
| 1.6 Architettura di una rete neurale | 21 |
| 2 L'apprendimento delle reti neurali | 25 |
| 2.1 Paradigmi ed addestramento | 25 |
| 3 Feed Forward | 29 |
| 4 Ricorrenti | 30 |
| Riferimenti bibliografici | 31 |

Introduzione

prova

Capitolo 1

Presentazione delle reti neurali artificiali

1.1 Cenni storici ed origini

Le reti neurali artificiali, dall'inglese *Artificial Neural Network (ANN)*, e la conseguente computazione neurale, sono state costruite ispirandosi ai sistemi neurali biologici, con l'obiettivo di modellarne il comportamento e la struttura e di simularne le funzioni basilari e fondamentali.

Una definizione semplice e formale di tali strutture è fornita dall'inventore del primo neurocomputer, il Dr. Robert Hecht-Nielsen che le definì come:

«...a computing system made up of a number of simple, highly interconnected processing elements, which process information by their dynamic state response to external inputs»

(da Neural Network Primer: Parte uno, Maureen Caudill, 1989)

Tradotto in italiano

«... un sistema di calcolo costituito da una serie di semplici, altamente interconnessi elementi di elaborazione, che elaborano le informazioni attraverso il

loro stato dinamico rispondendo agli input esterni.»

Le basi per lo studio di tali reti furono poste dallo psichiatra Warren McCulloch e del matematico Walter Pitts, i quali riuscirono a riprodurre una rete neurale utilizzando semplici circuiti elettrici collegati tra loro. Questa collaborazione portò alla luce l'analogia che sussiste tra le reti neurali e la macchina di Turing ed in tal modo si capì che qualsiasi operazione eseguita da una rete neurale poteva essere eseguita anche da un computer; infatti le reti che furono prodotte risultavano essere automi a stati finiti in grado di realizzare la logica delle proposizioni e di formulare ipotesi sulla natura dei meccanismi cerebrali, il tutto equivalentemente ad un programma per computer. Il frutto di tale lavoro fu reso noto nel 1943 con la pubblicazione del libro *"A logical calculus of the ideas immanent in nervous activity"*, nel quale fu schematizzato un combinatore lineare a soglia con dati binari multipli in entrata e un singolo dato binario in uscita.

Ciò che quindi McCulloch e Pitts riuscirono a fare è presentare un modello di neurone formale dimostrando che reti formate da tali neuroni riuscivano a computare funzioni della logica del primo ordine.

Un punto di svolta per lo studio delle reti neurali si ebbe successivamente alla pubblicazione del lavoro dello psicologo Donald Hebb, *"The Organization of Behavior"* nel 1949. Hebb trovò una correlazione tra la psicologia del comportamento umano e la fisiologia del sistema nervoso, donando un grande contributo alla teoria sull'apprendimento associativo, teoria che risultò essere alla base dei metodi di apprendimento delle reti neurali, e che si basava sulla nota legge di Hebb: *«Se un neurone A è abbastanza vicino ad un neurone B da contribuire ripetutamente e in maniera duratura alla sua eccitazione, allora ha luogo in entrambi i neuroni un processo di crescita o di cambiamento metabolico tale per cui l'efficacia di A nell'eccitare B viene accresciuta»*.

Il decennio che va dagli anni cinquanta agli anni sessanta fu totalmente influenzato dalla legge di Hebb a tal punto che numerosi gruppi di ricerca condussero esperimenti e test sulle funzionalità del cervello, fino a porre le basi per la nascita dell'intelligenza artificiale (AI).

CAPITOLO 1. PRESENTAZIONE DELLE RETI NEURALI ARTIFICIALI 8

Nel 1958 Frank Rosenblatt introdusse il primo schema di rete neurale che designò con il termine “**perceptron**”, in italiano percettrone, allo scopo di fornire un’interpretazione dell’organizzazione generale dei sistemi biologici attraverso un modello mirato all’analisi, in forma matematica, di funzioni quali ad esempio l’immagazzinamento delle informazioni. Il percettrone fu il primo modello di apprendimento supervisionato e presupponeva uno strato di ingresso ed uno di uscita, discriminando gli ingressi in due insiemi linearmente separabili e basandosi su una regola di apprendimento che si appellava alla minimizzazione dell’errore.

Il percettrone ancora oggi viene utilizzato in varie applicazioni e risultò essere un modello più efficace rispetto a quello binario di McCulloch e Pitts, poichè i suoi pesi sinaptici sono variabili e quindi in grado di apprendere.

Nonostante, però, l’iniziale successo di tale modello e l’interesse mostrato dalla comunità scientifica, tale rete neurale non risultò abbastanza potente; le reti a due strati basate sui percettroni avevano limiti operativi, infatti non riuscivano a risolvere tutte le classi di problemi, in particolare quelli non caratterizzati dalla separabilità lineare delle soluzioni come ad esempio l’operatore *XOR*, ovvero la funzione *OR* esclusivo, che discrimina gli ingressi in modo non linearmente separabile.

Solo negli anni ottanta, con il matematico Paul Werbos, si superarono i limiti del percettrone di Rosenblatt. Quest’ultimo introdusse uno o più livelli intermedi all’interno delle reti neurali creando una classe chiamata “**Multi-Layers Perceptron**”, ovvero percettrone multistrato, il cui metodo di addestramento principale era l’*error backpropagation*, ovvero l’algoritmo di retropropagazione dell’errore, che permetteva la modifica sistematica dei pesi delle connessioni, in modo da rendere la risposta della rete quanto più vicina a quella desiderata.

Tale algoritmo, proposto nel 1986 da David E. Rumelhart, G. Hinton e R. J. Williams consentì di superare le problematiche legate al percettrone di Rosenblatt e permise di risolvere il problema della separabilità non lineare delle soluzioni, rendendo quindi possibile calcolare la funzione *XOR* e segnando il definitivo rilancio delle reti neurali.

1.2 Analogie con il sistema nervoso ed applicazioni delle ANN

Come accennato nel paragrafo precedente, una rete neurale è un sistema computazionale costruito basandosi sui processi biologici naturali, il cui obiettivo è la riproduzione delle attività tipiche del cervello umano, ad esempio la comprensione del linguaggio, la percezione di immagini, il riconoscimento di forme

In altre parole, lo scopo di una rete neurale artificiale è l'emulazione del sistema nervoso animale, in particolar modo di quello umano, il quale presenta numerose caratteristiche che risultano essere ottime per la riproduzione di un sistema computazionale che lo imiti: è flessibile poiché si adatta ad ogni tipologia di situazione imparando, è robusto in quanto le cellule nervose muoiono ogni giorno senza avere effetti significativi sulla performance del sistema, è resistente, piccolo e dissipa poca energia.

Tale riproduzione non deve essere intesa come atto alla costruzione di un cervello artificiale. Infatti le caratteristiche delle reti neurali biologiche, riprese dalla computazione neurale artificiale, sono in esigua minoranza: gli stessi neuroni artificiali sono solo un'approssimazione dei neuroni biologici e sono in grado di riprodurre solo tre dei circa centocinquanta processi che sono tipici dei neuroni del cervello umano.

Il sistema nervoso umano può essere pensato, quindi, come una grande struttura computazionale formata da milioni di unità fortemente interconnesse tra loro in modo parallelo e che riesce a trasformare continui input in output ragionevoli. Le reti neurali rappresentano una riproduzione significativa di tale struttura, in particolar modo degli algoritmi di apprendimento e di ottimizzazione, basati su un modello *connessionistico* di calcolo: le operazioni responsabili dello scambio di informazioni avvengono per mezzo dell'interazione tra le unità elementari. Sono sistemi altamente paralleli: fornendo i dati del problema alle unità di input, la computazione si propaga in parallelo nella rete fino alle unità di output che producono il risultato.

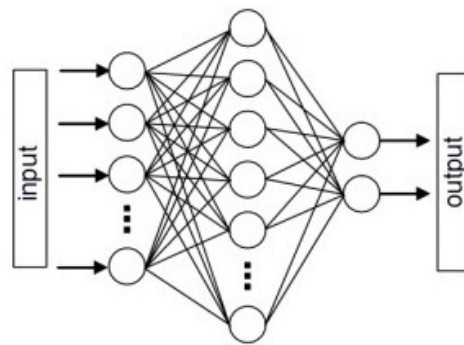


Figura 1.1: Esempio di rete neurale

Il campo tipico di applicazione delle reti neurali è quello dell'individuazione di legami di **ingresso-uscita** all'interno di sistemi complessi in cui i dati a disposizione sono molto numerosi o contenenti poca informazione e quindi non risulta chiaro quali relazioni deterministiche esistano tra le diverse variabili che caratterizzano il problema. Si tratta di campi in cui l'analisi statistica di tutte le variabili risulta difficoltosa o dispendiosa in termini di calcolo.

Ecco un breve elenco di applicazioni note:

- Elaborazione delle immagini (compressione e miglioramento di immagini in tempo reale);
- Elaborazione del suono (compressione e miglioramento in tempo reale della voce e della musica, aiuti audio e protesi ecc...);
- Rivelazione, verifica e identificazione di caratteri e oggetti (codici a barre, impronte digitali, firme, facce ecc...);
- Riconoscimento della voce (elaborazione di parole dettate, identificazione di persone per via vocale, selezione telefonica ecc...);
- Applicazioni basate su ingressi di origine sensoriale (ingresso tattile, visivo, acustico, olfattivo);
- Archivi (musica, video, immagini ecc...);
- Sintesi e miglioramento dati (fax, grafica, cancellazione di rumore, apparecchiature televisive ad alta risoluzione ecc...);

- Robot autonomi.

1.3 Il neurone biologico

Per convincersi a fondo sulla connessione tra le reti neurali artificiali e le reti neurali del sistema nervoso umano, e le numerose analogie che ne derivano, è bene dare una breve descrizione del secondo sistema, in modo da far emergere le principali caratteristiche ed i principali costituenti che sono utili al fine della comprensione delle reti neurali.

Il sistema nervoso è diviso in tre stadi:

- I *recettori*, che convertono gli stimoli esterni in impulsi elettrici che vengono inviati alla rete neurale;
- La *rete neurale*, la quale riceve gli impulsi elettrici provenienti dai recettori e li immagazzina come informazioni utili per poter prendere delle decisioni che vengono inviate, anch'esse sotto forma di impulsi elettrici, agli attuatori;
- Gli *attuatori* responsabili della trasformazione degli impulsi in risposte per l'ambiente esterno.

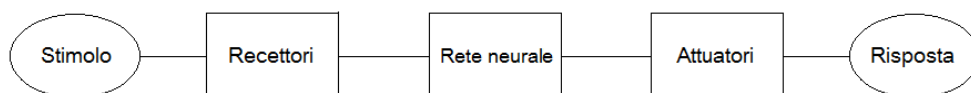


Figura 1.2: Stadi del sistema nervoso

Nel sistema nervoso centrale sono presenti circa 10^{11} cellule nervose, i *neuroni*. Tali cellule sono connesse strettamente tra loro attraverso numerosi collegamenti che nel loro insieme vanno a formare la rete neurale.

Un neurone è formato da un corpo centrale, identificato con il nome di

soma, all'interno del quale è presente il nucleo, e da molti prolungamenti citoplasmatici, detti *neuriti*, che si distinguono in:

- *dendriti*, organizzati con diramazioni ad albero che costituiscono il *ramo dendritico*;
- *assone*, la cui parte finale prende il nome di *bottone sinaptico*.

I dendriti ricevono segnali dai neuroni afferenti e li propagano verso il nucleo. L'assone conduce, invece, il segnale verso altre cellule grazie alla presenza del bottone sinaptico alla sua estremità, che risulta essere un'ulteriore ramificazione: quest'ultima va a formare i terminali attraverso i quali i segnali elettrici vengono trasmessi.

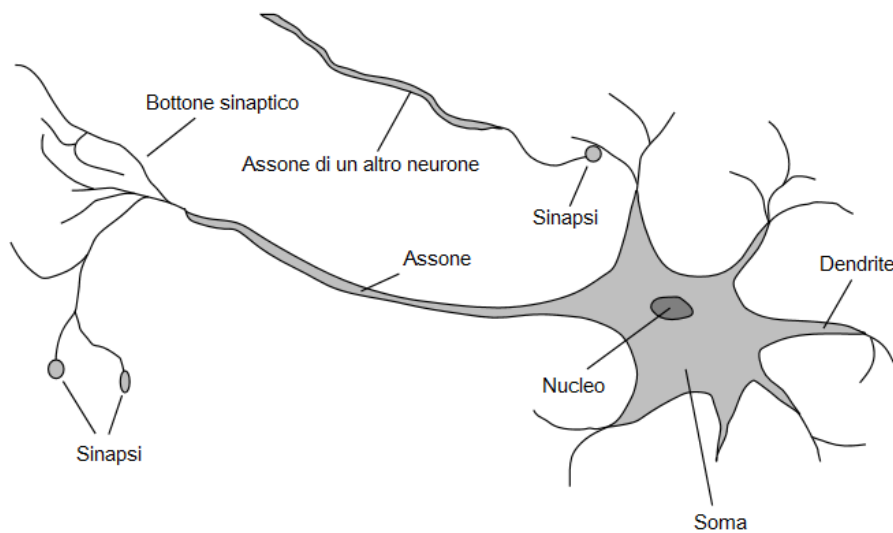


Figura 1.3: Neurone biologico

Tra un terminale di un assone e la cellula ricevente, esiste uno spazio che viene superato dagli impulsi elettrici per mezzo dei *neurotrasmettitori*, sostanze chimiche che guidano le informazioni tra le cellule e che si suddividono in *eccitatori* se favoriscono la creazione dell'impulso e *soppressori* se inibiscono l'impulso.

Il punto di connessione tra il terminale di un neurone ed il dendrite di un altro costituisce una struttura altamente specializzata che prende il nome di *sinapsi* e che risulta quindi essere la responsabile delle interazioni. La sinapsi svolge due processi: un processo *presinaptico* in cui viene liberato il neurotrasmettitore ed un processo *postsinaptico*, azionato dal neurotrasmettitore, che rigenera il segnale elettrico. Quindi una sinapsi converte un segnale elettrico in chimico durante il processo postsinaptico per poi riconvertirlo in elettrico durante la postsinapsi.

Un neurone trasmette un impulso elettrico lungo il suo assone nel momento in cui si verifica una differenza di potenziale elettrico tra l'interno e l'esterno della cellula, provocando la liberazione di un neurotrasmettitore di cui sopra.

1.4 Modello del neurone artificiale

La breve illustrazione del sistema neurale umano nel paragrafo precedente suggerisce uno schema per delineare l'organizzazione delle reti neurali artificiali: queste ultime sono formate da un elevato numero di unità computazionali, che possono essere equiparate ai neuroni umani, capaci di eseguire una somma pesata. Tali unità sono collegate tra loro attraverso delle connessioni, così come le sinapsi collegano i neuroni nella rete umana.

Consideriamo una generica unità j costituita da n canali di ingresso x_1, x_2, \dots, x_n . Gli input provenienti da strati precedenti o direttamente dall'esterno entrano nel neurone tramite tali canali.

Sulle connessioni sono presenti dei *pesi sinaptici* w_i , numeri reali che denotano l'*efficacia sinaptica*, ovvero la forza della connessione. Se $w_i > 0$ il canale è detto *eccitatorio*, se $w_i < 0$ il canale è *inibitorio*.

I segnali in entrata, pesati dalle rispettive sinapsi, sono convogliati nel *soma* del neurone artificiale, all'interno del quale vengono sommati producendo una combinazione lineare così definita:

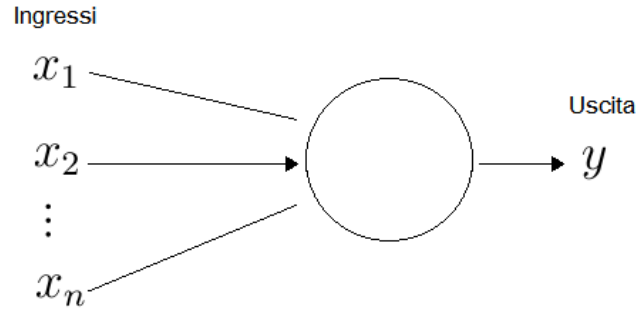


Figura 1.4: Canali di ingresso in un neurone

$$\sum_{i=1}^n w_i x_i \quad (1.1)$$

La somma pesata degli ingressi viene indicata con la parola *net* ed il segnale con cui il neurone trasmette la sua attività all'esterno è calcolato applicando una *funzione di attivazione* φ che limita l'ampiezza dell'output; si assume per comodità che le ampiezze degli output appartengono all'intervallo $[0, 1]$ oppure $[-1, 1]$.

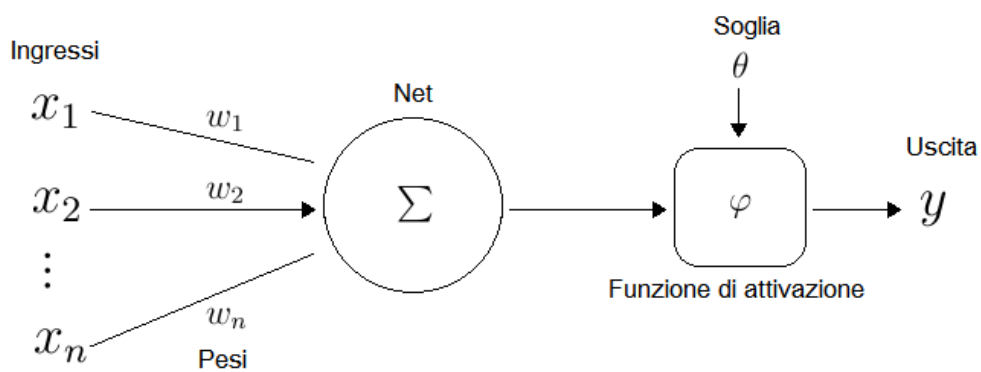


Figura 1.5: Modello neurone

Il modello neuronale include anche un valore *soglia* che ha l'effetto, a seconda della sua positività o negatività, di aumentare o diminuire il valore in ingresso alla funzione di attivazione.

L'output finale sarà allora:

$$y = \varphi\left(\sum_{i=1}^n w_i x_i\right) \quad (1.2)$$

E se indichiamo con θ il valore di soglia, la (1.2) diventerà:

$$y = \varphi\left(\sum_{i=1}^n w_i x_i - \theta\right) \quad (1.3)$$

Interpretando la soglia come il peso associato ad un ulteriore canale di ingresso x_0 , e quindi $w_0 = \theta$, potremmo anche scrivere :

$$y = \varphi\left(\sum_{i=0}^n w_i x_i\right) \quad (1.4)$$

Il modello finale sarà:

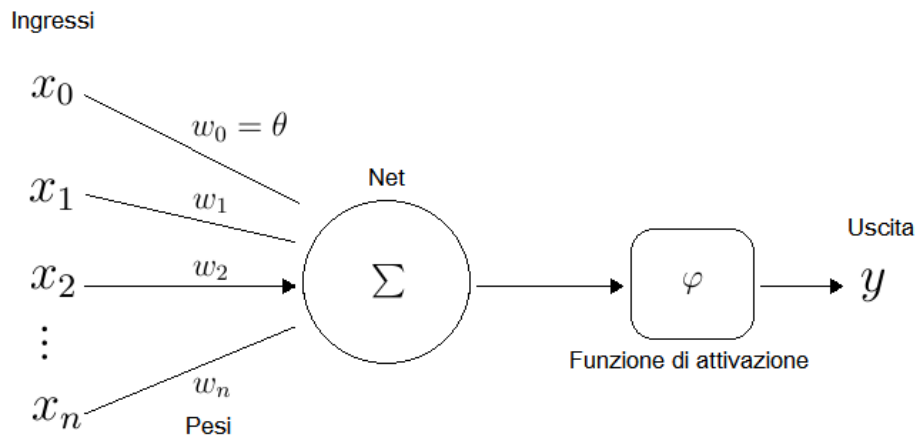


Figura 1.6: Modello neurone con soglia

L'effetto di un segnale x_i sul neurone è quindi uguale al prodotto $w_i \cdot x_i$ dove w_i è il peso attribuito alla sinapsi corrispondente ed il potenziale di attivazione è dato dalla somma algebrica dei prodotti di tutti i segnali di ingresso e dei valori dei pesi corrispondenti.

Schematizzando, ed individuando il neurone artificiale come l'unità di calcolo fondamentale della rete neurale, gli elementi base che lo rappresentano sono:

- Un insieme di connessioni;
- Un sommatore;
- Una funzione di attivazione;
- Un valore di soglia.

1.5 Funzioni di attivazione

La funzione di attivazione determina il tipo di risposta che un neurone è in grado di emettere. Definisce, quindi, l'uscita di un neurone in funzione del livello di attivazione. L'uscita può essere un numero reale, un numero reale appartenente ad un intervallo, oppure un numero appartenente ad un insieme discreto.

1.5.1 Funzione a soglia

Imponendo $a = \sum_{i=0}^n w_i x_i$, il valore di uscita di un neurone assunto tramite una funzione a soglia é:

$$y = \varphi(a) = \begin{cases} 1 & \text{se } a \geq 0 \\ 0 & \text{se } a < 0 \end{cases} \quad (1.5)$$

ed il grafico relativo

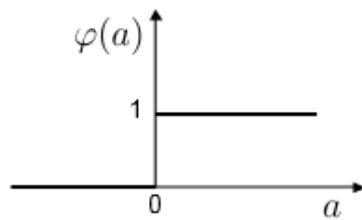


Figura 1.7: Funzione di attivazione a soglia

Se si traslascia il contributo dell'ingresso x_0 dal livello di attivazione e si ha quindi $a = \sum_{i=1}^n w_i x_i$, si avrà

$$y = \varphi(a) = \begin{cases} 1 & \text{se } a \geq \theta \\ 0 & \text{se } a < \theta \end{cases} \quad (1.6)$$

Il rispettivo grafico sarà

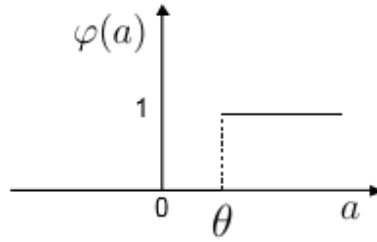


Figura 1.8: Funzione a soglia senza contributo x_0

A volte è opportuno che la funzione di attivazione assuma valori tra -1 e $+1$ ed in questo caso la funzione a soglia, che diviene la ben nota *funzione segno* è ridefinita così :

$$y = \varphi(a) = \begin{cases} 1 & \text{se } a > 0 \\ 0 & \text{se } a = 0 \\ -1 & \text{se } a < 0 \end{cases} \quad (1.7)$$

1.5.2 Funzione lineare

Se $a = \sum_{i=0}^n w_i x_i$, si avrà:

$$y = \varphi(a) = a \tag{1.8}$$

Il grafico rispettivo

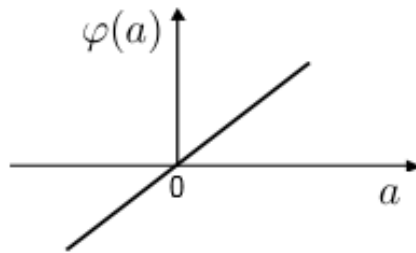


Figura 1.9: Funzione lineare

1.5.3 Funzione lineare a tratti

Un esempio di funzione di attivazione lineare a tratti è:

$$y = \varphi(a) = \begin{cases} 1 & \text{se } a \leq -0.5 \\ a + 0.5 & \text{se } -0.5 < a < 0.5 \\ 0 & \text{se } a \geq 0.5 \end{cases}$$

(1.9)

Che rappresentata è

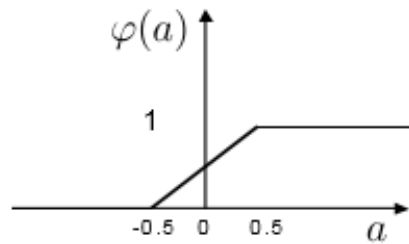


Figura 1.10: Funzione a tratti

1.5.4 Funzione sigmoide

La funzione sigmoide appartenente alla famiglia di funzioni continue non lineari ed è tra le più utilizzate. Tra queste funzioni riportiamo la *funzione logistica* così definita:

$$y = \varphi(a) = \frac{1}{1 + e^{-a}}$$

(1.10)

imponendo sempre $a = \sum_{i=0}^n w_i x_i$.

Essa assume questa forma

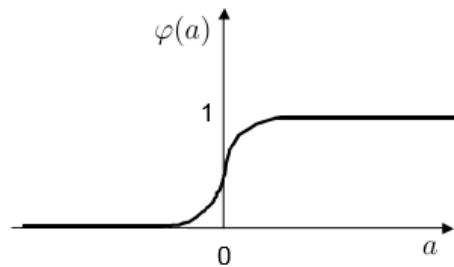


Figura 1.11: Funzione sigmoide

1.6 Architettura di una rete neurale

L'architettura di una rete neurale artificiale è caratterizzata da:

- Numero di strati di sinapsi;
- Numero di neuroni presenti nell'*input layer*, ovvero lo strato di ingresso;
- Numero di neuroni nell'*output layer*, lo strato di uscita.

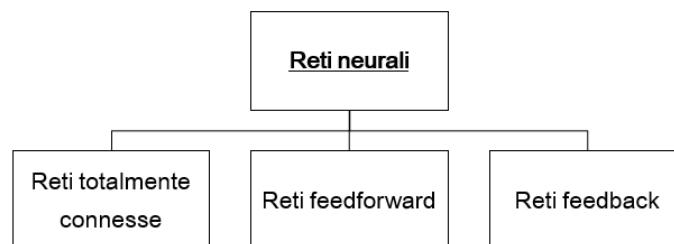


Figura 1.12: Suddivisione delle reti neurali

Le reti neurali si suddividono principalmente in due grandi classi: le reti *feedforward* e le reti *feedback* o *ricorrenti*. Si ha inoltre un'altra tipologia di reti che sono le *reti completamente connesse*.

Reti completamente connesse

Nelle reti completamente connesse ogni neurone è connesso con tutti gli altri. Le connessioni tra i neuroni sono bidirezionali e possono essere rappresentate per mezzo di una matrice quadrata W , di dimensione pari al numero di neuroni. Un suo generico elemento $w_{i,j}$ rappresenta il peso della connessione tra il neurone i ed il neurone j .

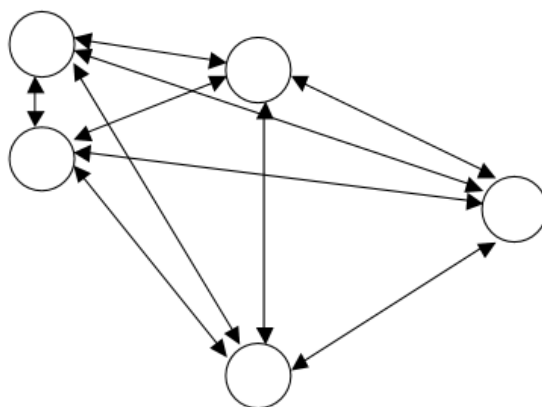


Figura 1.13: Rete completamente connessa

Reti feedforward a strati

In queste reti i segnali viaggiano dallo strato di ingresso verso lo strato di uscita e pertanto vengono anche chiamate **reti feedforward**. Nelle reti stratificate non esistono connessioni tra i neuroni all'interno di uno stesso strato, nè tra neuroni di strati non limitrofi. Ogni neurone in un generico strato è connesso con tutti quelli dello strato successivo ed i neuroni dello strato di ingresso hanno come unico compito quello di trasmettere i segnali ricevuti allo strato successivo, all'interno di essi non avviene alcuna computazione. Le reti feedforward a strati si distinguono in base al numero di strati che presentano, numero che dipende dallo specifico problema che si intende ri-

solvere.

- *Reti feedforward ad uno strato* : questa è una forma semplice di reti a strati. Il segnale nella rete si propaga in avanti senza cicli, ci sono connessioni che tornano indietro e nemmeno connessioni trasversali nel layer di output.

- *Reti feedforward a più strati* : le reti a più strati sono anche dette *reti multilivello* , in inglese **Multi-Layer Perceptron, MLP**. Tra l'input layer e l'output layer sono presenti uno o più strati di neuroni nascosti, si parla quindi di *hidden layers*. Nelle MLP non esistono connessioni nè tra neuroni di uno stesso strato nè tra neuroni di strati non limitrofi. Ogni strato ha connessioni entranti dal precedente strato e uscenti in quello successivo, quindi la propagazione del segnale avviene in avanti in modo aciclico e senza connessioni trasversali.

Tali reti vengono utilizzate per superare problemi che possono sorgere nella discriminazione dei segnali: attraverso i neuroni nascosti si ottengono delle rappresentazioni interne dei segnali di input che consentono il riconoscimento di forme più complesse, facilitando il compito della rete. Nonostante tutto, però, l'aggiunta di ulteriori strati nascosti non ottimizza le abilità di discriminazione della rete; questo è possibile solo se la funzione di attivazione è non lineare (una rete multistrato a neuroni lineari è sempre riconducibile ad una rete con due soli strati).

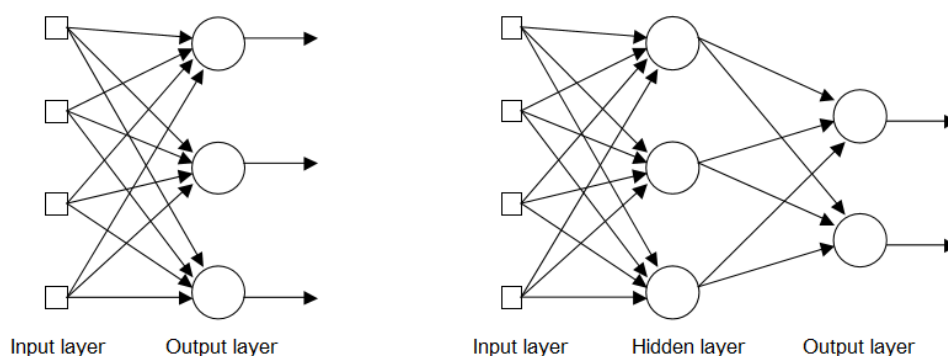


Figura 1.14: Rete ad uno strato (a sinistra), rete multilivello (a destra)

Reti ricorrenti

Una rete neurale ricorrente (RNN) si distingue dalle precedenti nel fatto che è ciclica. Il sistema è dinamico, le unità di output sono connesse con quelle intermedie e con quelle di input, si ha una *retroazione*. Dato un determinato stimolo, la risposta della rete ricorrente non viene dettata soltanto dai caratteri strutturali della rete stessa, come si verifica nella rete feedforward, ma varia in funzione del precedente contesto in cui si è manifestato lo stimolo. L'output non è determinato solo dall'input, ma anche da una *cronologia* di input che fornisce una forma di memoria a breve termine. Queste reti sono adatte per dati strutturati temporalmente ed hanno anche la possibilità di avere un'attività in assenza di una stimolazione input, poiché gli impulsi che viaggiano sulle vie ricorrenti possono essere sufficienti a mantenere il sistema in attività.

Capitolo 2

L'apprendimento delle reti neurali

2.1 Paradigmi ed addestramento

Per costruire una rete neurale efficiente, un passo fondamentale è individuare un *insieme di apprendimento* ed un *algoritmo di apprendimento*.

Per insieme di apprendimento si intende una collezione di esempi dai quali la rete può attingere per raggiungere lo scopo per il quale è stata progettata; per algoritmo, invece, un procedimento che permetta di prelevare le informazioni dall'insieme di apprendimento e di fissare dei parametri che vengano poi modificati attraverso operazioni iterative interfaccianti con l'ambiente. Le reti neurali si ispirano al tratto caratteristico del sistema nervoso ovvero la capacità di acquisire esperienza da esempi del mondo reale: per questo nella fase di apprendimento si parla di *addestramento* delle reti neurali attraverso dei *paradigmi di apprendimento*.

I paradigmi di apprendimento si suddividono in:

- *apprendimento supervisionato*
- *apprendimento non supervisionato*

Nell'apprendimento supervisionato si utilizza un *training set* cioè un set di esempi nel quale sono presenti coppie del tipo (x_k, y_{dk}) dove la prima variabile indica il k -esimo ingresso e la seconda la k -esima uscita. Con y_k si indica l'uscita reale e la si confronta con l'uscita desiderata: l'obiettivo è modificare i pesi affinché si minimizzi la differenza tra le due uscite. Il training set iniziale viene proposto ripetutamente finché $y_{dk} \approx y_k$, ovvero l'uscita desiderata sia il più simile possibile a quella reale, il tutto modificando i pesi in base alla legge di apprendimento scelta.

FIGURA

Un esempio di utilizzo di apprendimento supervisionato è quello adoperato nelle reti neurali che vengono utilizzate come *classificatori*, in grado cioè di riconoscere oggetti appartenenti a diverse categorie. Tale funzionamento si basa su vettori di numeri reali forniti all'ingresso che la rete deve poi attribuire ad una specifica classe. Un'altra classe di problemi risolti dalle reti neurali con apprendimento supervisionato è quella dell'approssimazione di funzioni. Data una funzione nota per punti, di cui non si conosce la forma analitica, il compito della rete neurale è quello di approssimarla nel modo migliore possibile, calcolandone il valore anche in punti diversi da quelli proposti in ingresso.

Nell'apprendimento non supervisionato la rete modifica i pesi autonomamente, si auto-organizza. Viene fornito solo il training set senza precisare le uscite.

GENERALIZZAZIONE

ASTRAZIONE

MEMORIA

PRINCIPIO DI OCCAM L'addestramento avviene sempre grazie ad un certo numero di esempi prelevati dal mondo reale. Nell'addestramento supervisionato, gli esempi forniti alla RN sono delle coppie (ingresso, uscita corrispon-

dente desiderata). . Altra classe di problemi risolti dalle RN è l'inseguimento (approssimazione) di funzioni. Una funzione in una o più variabili è nota per punti (non ne conosciamo la forma analitica); la RN deve riuscire ad approssimarla nel modo migliore possibile, calcolandone il valore anche in punti diversi da quelli proposti in ingresso. Nella fase di addestramento, diremo alla RN quale valore di uscita corrisponde a ciascun ingresso assegnato. Quanto nella modalità operativa riceverà un ingresso diverso, la rete dovrà tirare fuori un valore approssimato dell'uscita quanto più vicino possibile a quello reale. Nell'addestramento non supervisionato forniamo alla rete dei soli valori di ingresso (il TS, training set) senza precisare le relative uscite. La rete è in grado di auto-organizzarsi (come nelle SOM cui si è accennato prima) modificando la propria conoscenza interna (ovvero la memoria locale dei neuroni). Fra questi due estremi c'è l'addestramento graded: tipicamente la rete evolve in maniera autonoma, come nella modalità non supervisionata. I particolari istanti di tempo viene però utilizzata anche l'informazione sull'uscita desiderata. Si parla anche di reinforcement: ad intervalli di cicli di addestramento scanditi da un periodo si affianca alla modalità non supervisionata una modalità supervisionata. Noi non considereremo questo genere di addestramento.

Il corretto funzionamento di una neurale sull'insieme di apprendimento non offre, ovviamente, garanzia di un altrettanto soddisfacente funzionamento su altri dati relativi allo stesso concetto, ma non utilizzati nella fase di apprendimento (insieme che può essere utilizzato per effettuare un test). Inoltre, è evidente che l'architettura della rete neurale gioca un ruolo fondamentale per l'efficienza della fase di apprendimento. Si consideri, ad esempio, il caso delle reti feedforward, che si vedranno più avanti. Verrà enunciata la loro capacità universale di approssimazione, che le porta ad apprendere qualsiasi funzione, si potrebbe dire che sono in grado di imparare ogni complesso legame ingresso-uscita. Quando il numero delle unità cresce aumenta il potere computazionale, ma la capacità di generalizzare su nuovi esempi tende a diminuire dato che il fitting, in questo caso può essere definito come overfitting, sull'insieme di apprendimento ha luogo in un enorme spazio di parametri vincolati solo da pochi esempi. Questo origina una sorta di principio

di indeterminazione dell'apprendimento secondo il quale non è possibile al variare dei pesi della rete neurale ottenere un'adeguata generalizzazione per nuovi esempi. Come si vedrà la limitazione del numero degli ingressi risulta particolarmente importante per limitare l'over-fitting. La memorizzazione e la generalizzazione possono anche essere espresse in termini di dispendio di risorse di memoria richiesto nel primo caso rispetto al secondo. Quindi l'efficienza nell'uso delle risorse disponibili, principio di base per un ingegnere, è anche il principio a cui si adegua il sistema nervoso nell'acquisizione di esperienza.

Tale principio è ben noto in letteratura tecnica sotto il nome di "Principio di Occam" e può essere enunciato nel seguente modo in relazione alle reti neurali:

"date due reti che soddisfano l'insieme di apprendimento, la rete di minore complessità è quella che si comporta meglio su esempi non visti, cioè ha la migliore capacità di generalizzazione"

Il Principio di Occam è la guida più adatta alla determinazione della struttura della rete neurale ottimale dal punto di vista della capacità di generalizzazione. Infatti, esso suggerisce, nel caso delle reti neurali supervisionate, di utilizzare una funzione obiettivo da minimizzare costituita da due termini: l'uno riguardante l'adeguatezza della rete rispetto all'insieme di apprendimento; l'altro riguardante la complessità della rete. Indicando con E_s il primo, con E_c il secondo e con m un valore che misura il numero dei parametri della rete, si ha per la funzione obiettivo:

Capitolo 3

Feed Forward

Capitolo 4

Ricorrenti

Bibliografia

- [1] Prof.ssa Lazzarini Beatrice, (2015), *Introduzione alle reti neurali*.
- [2] Prof. Gambosi Giorgio, (2010), *Reti neurali, note dal corso di Machine Learning*.
- [3] Prof.ssa Labonia Laura, *Storia delle reti neurali artificiali*.
- [4] Prof. Bicego Manuele, *Riconoscimento e recupero dell'informazione per bioinformatica. Reti neurali*.
- [5] Ing. Pioggia Giovanni, (2009), *Modelli di sistemi fisiologici*