

## Reti neurali e deep learning

Le reti neurali sono una famiglia di predittori che, nella loro forma più semplice (detta rete neurale *feedforward*), sono ottenuti combinando semplici predittori della forma  $g(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x})$ . La funzione non lineare  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  è detta funzione di attivazione e tipicamente è definita come

$$\sigma(z) = \frac{1 - e^{-z}}{1 + e^{-z}} \in [-1, +1]. \quad (1)$$

Una rete neurale feedforward calcola una funzione  $f : \mathbb{R}^d \rightarrow \mathbb{R}^n$  rappresentata come un grafo diretto aciclico  $G = (V, E)$ . Un esempio particolare è quello mostrato in Figura 1, dove i nodi sono suddivisi in strati, c'è un solo nodo nello strato di output (*output layer*) e, ordinando gli strati da quello di input (*input layer*) a quello di output, ogni nodo non di input è connesso a tutti i nodi dello strato precedente. In generale, possiamo partizionare i nodi in tre sottoinsiemi:  $V = V_{\text{in}} \cup V_{\text{hid}} \cup V_{\text{out}}$ , dove  $V_{\text{in}}$  (con  $|V_{\text{in}}| = d$ ) sono i nodi dello strato di input (ovvero, nodi che non ricevono nessun arco in ingresso),  $V_{\text{out}}$  (con  $|V_{\text{out}}| = n$ ) sono i nodi nello strato di output (ovvero, nodi che non emettono nessun arco in uscita) e  $V_{\text{hid}}$  sono i nodi dello strato nascosto (*hidden layer*).

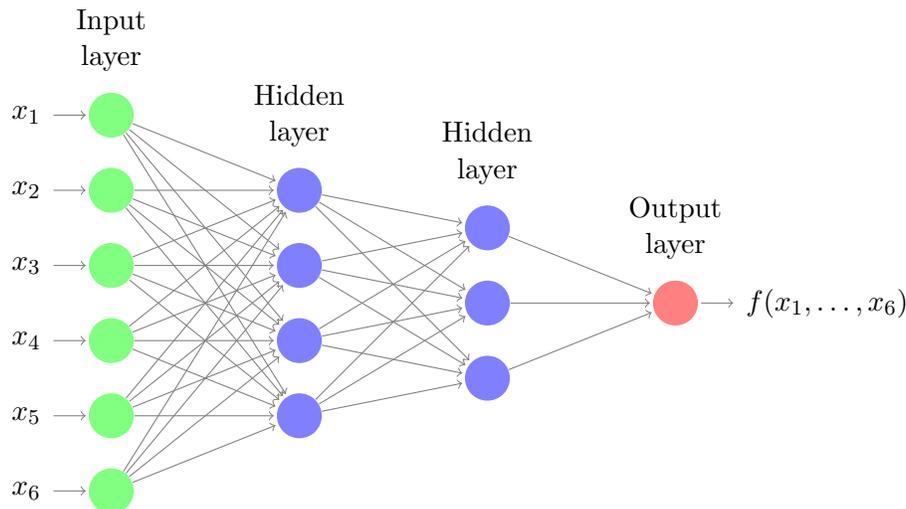


Figura 1: Un esempio di rete neurale feedforward con output singolo e due strati nascosti in cui le connessioni sono complete fra ogni coppia di strati adiacenti (per esempio fra i due strati nascosti) e assenti fra ogni coppia di strati non adiacenti (per esempio fra strato di input e strato di output).

Ad ogni arco  $(i, j) \in E$  è associato un parametro (detto peso)  $w_{i,j} \in \mathbb{R}$ ; indichiamo con  $W$  la matrice  $|V| \times |V|$  dei pesi, dove  $w_{i,j} = 0$  se  $(i, j) \notin E$ . Il grafo  $G$ , la matrice dei pesi  $W$  e la

funzione di attivazione  $\sigma$  definiscono la funzione  $\mathbf{f} = \mathbf{f}_{G,W,\sigma}$  calcolata dalla rete. Per calcolare questa funzione, ad ogni nodo  $i \in V$  è associato un valore  $v_i \in \mathbb{R}$  nel modo spiegato in seguito.

Dato  $j \in V \setminus V_{\text{in}}$ , indichiamo con:

- $\mathbf{w}(j)$  il vettore le cui componenti sono i pesi  $w_{i,j}$  per ogni  $(i,j) \in E$ ,
- $\mathbf{v}(j)$  il vettore le cui componenti sono i valori  $v_i$  per ogni  $(i,j) \in E$ .

Il valore  $\mathbf{f}(\mathbf{x}) = (f_1, \dots, f_n)$  viene calcolato nel modo seguente:

1. ogni nodo  $i \in V_{\text{in}}$  assume valore  $v_i = x_i$ ,
2. ogni nodo  $j \in V \setminus V_{\text{in}}$  assume valore  $v_j = \sigma(\mathbf{w}(j)^\top \mathbf{v}(j))$ ,
3.  $f_k = v_k$ , dove  $k$  è il  $k$ -esimo nodo di output.

Fissati  $d$  (dimensionalità dell'input) e  $n$  (dimensionalità dell'output), introduciamo la classe  $\mathcal{F}_{G,\sigma}$  dei predittori  $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^n$  tali che  $\mathbf{f} = \mathbf{f}_{G,W,\sigma}$  per una qualche matrice dei pesi  $W$ , ovvero predittori rappresentabili da una rete neuronale feedforward con grafo  $G$ , matrice dei pesi  $W$  e funzione di attivazione  $\sigma$ .

Chiaramente, possiamo utilizzare una rete neuronale con un singolo nodo di output come classificatore binario assegnando a un'istanza  $\mathbf{x}$  l'etichetta  $\text{sgn}(f_{G,W,\sigma}(\mathbf{x}))$ . Ci chiediamo allora quali classificatori binari  $f : \mathbb{R}^d \rightarrow \{-1, +1\}$  può realizzare una rete neuronale. Dimostriamo ora un risultato piuttosto sorprendente: quando le istanze sono binarie ( $\mathbf{x} \in \{-1, +1\}^d$ ) e la funzione di attivazione  $\sigma$  è la funzione segno  $\text{sgn}$ , una singola rete con un unico strato nascosto è in grado di implementare qualsiasi classificatore  $f : \{-1, +1\}^d \rightarrow \{-1, +1\}$ . Si noti che la restrizione a istanze binarie non è vincolante in pratica, in quanto ogni coordinata reale  $x_i$  è comunque rappresentata in un computer da una stringa di bit di lunghezza fissa.

**Teorema 1** *Per ogni  $d$ , esiste una rete neuronale  $G = (V, E)$  con  $d + 1$  nodi di input, uno strato di nodi nascosti e un nodo di output, tale che  $\mathcal{F}_{G,\text{sgn}}$  contiene tutte le funzioni della forma  $f : \{-1, +1\}^d \rightarrow \{-1, +1\}$ .*

**DIMOSTRAZIONE.** Fissiamo una rete con  $2^d + 1$  nodi nascosti e connettività completa dello strato di input con lo strato nascosto e dello strato nascosto con lo strato di output. Numeriamo i nodi come segue: 0 è il nodo di output,  $1, \dots, 2^d + 1$  sono i  $2^d + 1$  nodi nascosti e i seguenti sono i  $d + 1$  nodi di input. Dato una qualunque funzione  $f : \{-1, +1\}^d \rightarrow \{-1, +1\}$  dimostriamo che esiste una matrice dei pesi  $W$  tale che  $f_{G,W,\text{sgn}} = f$ . Siano  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , con  $N \leq 2^d$ , tutte e sole le istanze  $\mathbf{x}_i \in \{-1, +1\}^d$  tali che  $f(\mathbf{x}_i) = 1$ . Si osservi che per ogni  $\mathbf{x} \in \{-1, +1\}^d$  e per ogni  $i = 1, \dots, N$ ,

$$\begin{aligned} \mathbf{x}^\top \mathbf{x}_i &\leq d - 2 && \text{se } \mathbf{x} \neq \mathbf{x}_i \\ \mathbf{x}^\top \mathbf{x}_i &= d && \text{altrimenti.} \end{aligned}$$

Quindi,  $g_i(\mathbf{x}) = \text{sgn}(\mathbf{x}^\top \mathbf{x}_i - d + 1) = 1$  se e solo se  $\mathbf{x} = \mathbf{x}_i$ . Ogni funzione  $g_i$  può essere implementata dal nodo nascosto  $i$  nel modo seguente: si trasforma ogni istanza  $\mathbf{x} = (x_1, \dots, x_d)$  nell'istanza  $\tilde{\mathbf{x}} = (\mathbf{x}, 1) = (x_1, \dots, x_d, 1)$  —questa è la tecnica standard che abbiamo usato anche per i classificatori lineari— e si assegnano i pesi  $\mathbf{w}(i) = (\mathbf{x}_i, -d + 1)$ . Questo implica che  $\text{sgn}(\mathbf{w}(i)^\top \tilde{\mathbf{x}}) = g_i(\mathbf{x})$ . Poi assegnamo  $\mathbf{w}(N + 1) = (0, \dots, 0, 1)$  in modo che  $\text{sgn}(\mathbf{w}(i)^\top \tilde{\mathbf{x}}) = 1$  per ogni  $\mathbf{x}$  (i rimanenti  $\mathbf{w}(i)$  per  $i = N + 2, \dots, 2^d + 1$  possono essere assegnati arbitrariamente). Infine,  $f$  può essere rappresentata

come  $f(\mathbf{x}) = g_1(\mathbf{x}) \vee \dots \vee g_N(\mathbf{x})$ , il che corrisponde ad assegnare i pesi

$$\mathbf{w}(0) = \left( \underbrace{1, \dots, 1}_{N \text{ volte}}, \underbrace{N-1, 0, \dots, 0}_{2^d - N \text{ volte}} \right)$$

dove 0 è il nodo di output. □

La dimostrazione del teorema precedente utilizza una rete con un singolo strato nascosto di dimensione almeno pari al numero massimo di istanze positive per le funzioni  $f$  da apprendere. Questo numero è tipicamente esponenziale nella dimensione  $d$ . Il risultato seguente (che non dimostreremo) stabilisce che se vogliamo una rete in grado di apprendere tutte le funzioni da  $\{-1, +1\}^d$  a  $\{-1, +1\}$ , allora tale rete deve avere un numero complessivo di nodi esponenziale in  $d$ .

**Teorema 2** *Per ogni intero  $d$  sia  $s(d)$  il più piccolo intero tale che esiste  $G = (V, E)$  con  $|V| = s(d)$  per cui  $\mathcal{F}_{G, \text{sgn}}$  contiene tutte le funzioni della forma  $f : \{-1, +1\}^d \rightarrow \{-1, +1\}$ . Allora  $|V| = \Omega(2^{d/3})$ . Un risultato analogo vale quando  $\text{sgn}$  è rimpiazzata dalla funzione sigmoideale (1).*

Le reti feedforward con un singolo strato nascosto non solo implementano qualsiasi funzione booleana, ma sono anche in grado di approssimare qualsiasi funzione  $f : [-1, +1]^d \rightarrow [-1, +1]$  che sia Lipschitziana (ovvero tale che esiste  $L < \infty$  che soddisfa  $|f(\mathbf{x}) - f(\mathbf{x}')| \leq L \|\mathbf{x} - \mathbf{x}'\|$  per ogni  $\mathbf{x}, \mathbf{x}' \in [-1, +1]^d$ ) quando la funzione di attivazione utilizzata è la sigmoideale (1). Ovvero, per ogni  $\varepsilon > 0$  e per ogni  $f : [-1, +1]^d \rightarrow [-1, +1]$  Lipschitziana esistono  $G$  (con un singolo strato nascosto) e  $W$  tali che  $|f_{G, W, \sigma}(\mathbf{x}) - f(\mathbf{x})| \leq \varepsilon$  per ogni  $\mathbf{x} \in [-1, +1]^d$ . Come nel caso delle funzioni booleane, il numero di nodi nello strato nascosto necessario per rappresentare la funzione  $f$  con un grado di approssimazione fissato  $\varepsilon > 0$  può essere esponenziale in  $d$ .

Abbiamo quindi visto che un singolo strato nascosto è sufficiente per implementare qualsiasi funzione booleana o approssimare qualsiasi funzione Lipschitziana a patto di essere disposti a utilizzare un numero esponenziale di nodi. Il Teorema 2 ci dice che se vogliamo poter imparare un numero elevato di funzioni con una rete  $G$ , allora  $G$  deve avere un numero esponenziale di nodi. Resta aperta la possibilità che utilizzando una rete con un grande numero di strati nascosti ciascuno contenente relativamente pochi nodi (*deep neural network*) sia possibile imparare una classe sufficientemente ampia e interessante di funzioni. Questo è un tema attivo di ricerca. I risultati fino ad oggi disponibili sono promettenti, anche se molti dettagli sono ancora da investigare.

In una rete neuronale profonda, il grafo  $G$  e la funzione di attivazione  $\sigma$  dovrebbero essere scelti in modo tale che  $\mathcal{F}_{G, \sigma}$  contenga funzioni che siano utili rispetto alla tipologia di dati che si vogliono apprendere (immagini, testi, eccetera). Per esempio, le reti neurali convolutive sono una famiglia di reti progettate specificatamente per risolvere problemi di classificazioni di immagini rappresentate come matrici di pixel. In generale, comunque, la questione di quale coppia  $G, \sigma$  sia più adeguata per risolvere un dato problema di apprendimento è anch'esso un tema attivo di ricerca.

Occupiamoci ora del problema di addestrare una rete neuronale feedforward con funzione di attivazione sigmoideale e  $n = 1$  (un singolo nodo di output). Le reti feedforward vengono tipicamente addestrate usando l'algoritmo di discesa del gradiente stocastico,

$$w_{i,j} \leftarrow w_{i,j} - \eta_t \frac{\partial \ell_{Z_t}(W)}{\partial w_{i,j}}$$

dove  $Z_t$  è l'indice di un esempio estratto a caso dal training set e  $\ell_t(W) = \ell(f_{G,W,\sigma}(\mathbf{x}_t), y_t)$  è una funzione di perdita tale che  $\ell(\cdot, y)$  è convessa (per esempio, la funzione di perdita quadratica o la funzione di perdita logistica). Spesso, per accelerare la convergenza, viene anche usata la versione *mini-batched* del gradiente stocastico,

$$w_{i,j} \leftarrow w_{i,j} - \eta_t \frac{1}{|S_t|} \sum_{s \in S_t} \frac{\partial \ell_s(W)}{\partial w_{i,j}}$$

dove  $S_t$  è un sottoinsieme di cardinalità fissata estratto a caso dal training set.

È importante osservare che  $\ell_t(W)$  è una funzione non convessa dei pesi  $W$  anche quando  $\ell(\cdot, y)$  è convessa per ogni  $y$ . Questo ha due conseguenze importanti:

1. Il problema di determinare i pesi  $W$  che minimizzano (anche in modo approssimato) l'errore sul training set è considerato computazionalmente intrattabile.
2. Il metodo della discesa del gradiente converge tipicamente a un minimo locale del training error.

In pratica, si osserva che i minimi locali del training error corrispondono in generale a funzioni con basso rischio. Questo potrebbe significare che molti dei minimi locali sono vicini al minimo globale (che in assenza di overfitting corrisponde al modello di rischio minimo). La ragione precisa delle buone prestazioni del gradiente stocastico sulle reti neurali resta però un problema di ricerca aperto.

L'applicazione della discesa del gradiente per l'addestramento di una rete neurale feedforward prende il nome di algoritmo di retropropagazione dell'errore. Questo algoritmo utilizza la regola della derivazione di funzioni composte,

$$\frac{df(g(x))}{dx} = \frac{df(g)}{dg} \frac{dg(x)}{dx}. \quad (2)$$

Per comodità, supponiamo che 0 sia l'indice del nodo di output. Quindi  $f_{G,W,\sigma}(\mathbf{x}_t) = v_0 = \sigma(s_0)$ , dove  $s_0 = \mathbf{w}(0)^\top \mathbf{v}(0)$ . Per un qualunque  $i$  nel primo strato nascosto, ovvero tale che  $(i, 0) \in E$ , applicando due volte la regola (2) abbiamo

$$\frac{\partial \ell_t(W)}{\partial w_{i,0}} = \frac{\partial \ell(v_0, y_t)}{\partial w_{i,0}} = \frac{\partial \ell(v_0, y_t)}{\partial v_0} \frac{\partial v_0}{\partial s_0} \frac{\partial s_0}{\partial w_{i,0}} = \ell'(v_0) \frac{d\sigma(s_0)}{ds_0} \frac{\partial s_0}{\partial w_{i,0}}.$$

La derivata  $\ell'(v_0) = \frac{\partial \ell(v_0, y_t)}{\partial v_0}$  è calcolabile direttamente a partire dalla definizione della funzione di perdita. Per esempio,

$$\ell(x, y) = (x - y)^2 \quad \text{implica} \quad \frac{\partial \ell(x, y_t)}{\partial x} = 2(x - y).$$

In modo simile, la derivata  $\frac{d\sigma(s_0)}{ds_0} = \sigma'(s_0)$  è calcolabile direttamente a partire dalla funzione di attivazione. Per esempio,

$$\sigma(z) = \frac{1 - e^{-z}}{1 + e^{-z}} \quad \text{implica} \quad \sigma'(z) = \frac{(1 - \sigma(z))^2}{2}.$$

Infine, ricordando che  $s_0 = \mathbf{w}(0)^\top \mathbf{v}(0)$ , abbiamo

$$\frac{\partial s_0}{\partial w_{i,0}} = v_i = \sigma(\mathbf{w}(i)^\top \mathbf{v}(i)) .$$

Possiamo quindi scrivere

$$\frac{\partial \ell_t(W)}{\partial w_{i,0}} = \ell'(v_0) \sigma'(s_0) v_i .$$

Calcolamo ora  $\frac{\partial \ell_t(W)}{\partial w_{i,j}}$  per i nodi  $i$  nel secondo strato nascosto. Si noti che

$$v_0 = \sigma(\mathbf{w}(0)^\top \mathbf{v}(0)) = \sigma \left( \sum_{j:(j,0) \in E} w_{j,0} \sigma(\mathbf{w}(j)^\top \mathbf{v}(j)) \right) = \sigma \left( \sum_{j:(j,0) \in E} w_{j,0} \sigma \left( \sum_{i:(i,j) \in E} w_{i,j} v_i \right) \right) .$$

Quindi, per ogni nodo  $i$  nel secondo strato nascosto, denotando  $s_i = \mathbf{w}(i)^\top \mathbf{v}(i)$  possiamo scrivere

$$\frac{\partial \ell_t(W)}{\partial w_{i,j}} = \underbrace{\frac{\partial \ell(v_0, y_t)}{\partial s_0}}_{\ell'(v_0) \sigma'(s_0)} \underbrace{\frac{\partial s_0}{\partial v_j}}_{w_{j,0}} \underbrace{\frac{\partial v_j}{\partial s_j}}_{\sigma'(s_j)} \underbrace{\frac{\partial s_j}{\partial w_{i,j}}}_{v_i} = \ell'(v_0) \sigma'(s_0) w_{j,0} \sigma'(s_j) v_i .$$

Possiamo riscrivere quanto fatto in modo più compatto sfruttando il fatto che per ogni nodo  $i$ ,

$$\frac{\partial \ell_t(W)}{\partial s_i} = \sum_{j:(i,j) \in E} \frac{\partial \ell_t(W)}{\partial s_j} \frac{\partial s_j}{\partial v_i} \frac{\partial v_i}{\partial s_i} = \sum_{j:(i,j) \in E} \frac{\partial \ell_t(W)}{\partial s_j} w_{i,j} \sigma'(s_i) .$$

Quindi, introducendo la definizione ricorsiva

$$\delta_i = \frac{\partial \ell_t(W)}{\partial s_i} = \begin{cases} \ell'(v_0) \sigma'(s_0) & \text{se } i = 0 \\ \sigma'(s_i) \sum_{j:(i,j) \in E} \delta_j w_{i,j} & \text{altrimenti} \end{cases}$$

possiamo infine scrivere la derivata parziale per un qualunque  $(i, j) \in E$  come

$$\frac{\partial \ell_t(W)}{\partial w_{i,j}} = \frac{\partial \ell_t(W)}{\partial s_j} \frac{\partial s_j}{\partial w_{i,j}} = \delta_j v_i .$$

Addestrare una deep network usando l'algoritmo di retropropagazione dell'errore non è immediato. Per facilitare una convergenza rapida ad un buon minimo locale si utilizzano un insieme di euristiche (per esempio, inizializzazione casuale dei pesi in  $W$ ). Il corretto utilizzo di queste euristiche non è ancora compreso a fondo e deve quindi essere appreso dall'utente attraverso l'esperienza.