

Deep Learning

Daniele Avolio

October 9, 2023

Contents

1	Introduzione	4
2	Deep Learning 101	4
2.1	Architetture e strumenti nel deep learning	4
2.2	Libri utili	4
2.3	Strumenti che useremo	5
2.4	Schema generale di un problema di deep learning	5
2.5	Perché si usa il termine "Tensore"?	5
2.6	AI vs DL	6
3	Introduzione alle Reti Neurali	7
3.1	Il modello di McCulloch-Pitts	7
3.2	Modello di Rosenblatt	7
3.3	Esempio con 2 neuroni	8
3.4	Rappresentazione le funzioni logiche	9
3.4.1	AND	9
3.4.2	Il problema dello XOR	10
3.5	Gli ingredienti di una rete neurale	11
3.5.1	Il grafo g	11
3.5.2	La funzione di loss	12
3.5.3	Funzione di loss per Regressione	13
3.5.4	Funzione di loss per classificazione	13
3.6	L'ottimizzatore o	13
3.6.1	Il metodo di discesa del gradiente	14
3.6.2	Il metodo di discesa del gradiente stocastico	15
3.6.3	Linee guida sul learning rate	16
4	Classificazione Binaria	18
4.1	Caricamento e gestione del Dataset	19
4.2	Definizione della Rete Neurale	20
4.3	Plotting del modello	22
4.4	Training del modello e valutazione	23
4.4.1	Validation Set	23
4.5	Prediction	27
4.6	Come risolvere problemi di accuracy bassa?	27
4.7	Early Stopping	28
5	Classificazione Multiclasse	28
5.1	Descrizione del dataset - Reuters	28
5.1.1	Come processiamo l'output?	28
5.1.2	Softmax activation function	29

6	Lab: Introduzione Python	31
6.1	MatPlotLib	31
6.1.1	Plots	31
6.1.2	Sub-plots	31
6.2	NumPy	32
7	Lab: Reti Neurali da zero	34
7.1	Introduzione	34
7.2	Esempio pratico	34
7.2.1	Il grafo	35
7.2.2	La funzione loss	35
7.2.3	L'ottimizzatore	36
7.2.4	Discesa del gradiente	36
7.2.5	Inizializzatore	36
7.2.6	Fix Point Procedure	36
7.3	Esempio da zero	36
7.3.1	La funzione Sigmoid	36
7.4	Tangente Iperbolica TanH	37
7.5	ReLu	38
7.6	Scalare i valori	39
7.7	Plottare la loss	40

■ 1 Introduzione

Info esame: Tecnicamente, per ora rimane un progetto e potrebbe essere di gruppo. Ancora non ci sono informazioni precisissime.

■ 2 Deep Learning 101

In questo corso affronteremo diverse tematiche, il che può sembrare assurdo se ci si pensa.

- Classificazione (binaria, cioè sì o no)
- Multi-class classification (non più binaria)
- Regressione (il guessing viene fatto su un valore numerico)
- Gestione di immagini e riconoscimento
- Serie numeriche (predizioni di mercato e trend)
- Classificazione di testi

■ 2.1 Architetture e strumenti nel deep learning

- Autoencoder
 - Tutti i possibili tipi
 - Qui si fa anche **Clustering** e **Anomaly detection**
- Architetture generative
 - Tutti i possibili tipi
- XAI: Explainable AI

■ 2.2 Libri utili

- "Deep Learning in Python"
- "Tensorflow tutorial"

■ 2.3 Strumenti che useremo

- Tensorflow
 - High-level più di altri
- Keras
 - High-level API basato su Tensorflow
 - Ci saranno cose che non possiamo fare con Keras perché è troppo ad alto livello

■ 2.4 Schema generale di un problema di deep learning

Abbiamo delle coppie $(x_0, y_0), (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ dove x_i è un vettore di features e y_i è un valore numerico (regressione) o una classe (classificazione).

$$y_i = f(x_i)$$

Non conosciamo la funzione f , quindi dobbiamo impararla.

$$y = \alpha x + \beta$$

Con una rete neurale puoi approssimare praticamente qualsiasi funzione.

Una rete neurale permette di collegare un input di dati a una funzione di output.

Abbiamo diversi tipi di reti neurali a seconda del tipo di problema che vogliamo risolvere. È importante essere in grado di selezionare l'architettura giusta per risolvere il problema.

Definiamo alcuni concetti che useremo:

- N : rete neurale
- w : valori dei pesi della rete neurale
- f : funzione di output della rete neurale

$$f \in N(w)$$

■ 2.5 Perché si usa il termine "Tensore"?

Un *tensore* non è altro che una matrice.

- 0D tensor: scalar
- 1D tensor: vector
- 2D tensor: matrix
- 3D tensor: tensor

■ 2.6 AI vs DL

- AI: è un ampio insieme di tecniche per risolvere problemi che richiedono "intelligenza".
 - Esempio: Stockfish, un programma che gioca a scacchi.
- DL: È un sottoinsieme di AI che si concentra sull'*astrazione*.
 - L'astrazione consiste nel fornire una funzione che traduce dati di input in dati di output senza conoscere la funzione stessa.
 - È un approccio induttivo: fornisci un input e ti aspetti un output, senza conoscere la funzione che li collega.
 - Questo è completamente diverso dall'AI basata sulla logica.

■ 3 Introduzione alle Reti Neurali

■ 3.1 Il modello di McCulloch-Pitts

Un modello pensato dai due tizi qui presenti,

- Warren McCulloch
- Walter Pitts

Era formato da:

- Un insieme di neuroni
- Un insieme di connessioni tra i neuroni
- Un insieme di pesi associati alle connessioni
- Una funzione di attivazione
- Una funzione di output

Quindi immaginiamo x_1, x_2, \dots, x_n che vengono dati come input, e quello che viene fuori è un valore di $y \in \{0, 1\}$. Questo è un modello di classificazione binaria.

In soldoni, in deep learning si usa un vettore di numeri per tirare fuori un altro vettore di numeri.

Nella maggior parte del tempo, però, non lavoriamo solamente con i dati.

- immagini
- Testo
- Audio
- ...

Non ci sono concetti di foto, video, immagini. L'unico concetto che esiste è quello dei numeri.

■ 3.2 Modello di Rosenblatt

Qui il modello è leggermente diverso. Gli elementi in questo modello sono i seguenti:

- Valori di input: x_i

- Funzione di attivazione: ϕ
- Pesi degli archi: w_i
- Bias: b

$$h(x|w, b) = h\left(\sum_{i=1}^l w_i \cdot x_i - b\right) = h\left(\sum_{i=1}^l w_i \cdot x_i\right) = \text{sign}(w^T x) \quad (1)$$

Nota: Quando il **bias** non viene specificato, allora si assume che sia 0.

I pesi w_i sono collegati archi che vanno da x_i al prossimo neurone. Sia il valore di input che il peso sono **numeri reali**. Non sono lo stesso valore, hanno solamente il formato di *reale* che è uguale tra loro. Ciò che viene fatto è solamente la somma della prodotto tra ogni peso w_i e x_i **meno** il bias.

La funzione di attivazione: Dipende. Ogni funzione che ha *2 stati* va bene per noi. Una funzione di attivazione può essere una qualsiasi che in un punto ha valore 1 e in un altro ha valore -1.

■ 3.3 Esempio con 2 neuroni

Immaginiamo di avere un piano cartesiano con una retta che interseca in 2 punti.

$$\text{sign}(w_1 \cdot x_1 + w_2 \cdot x_2) \quad (2)$$

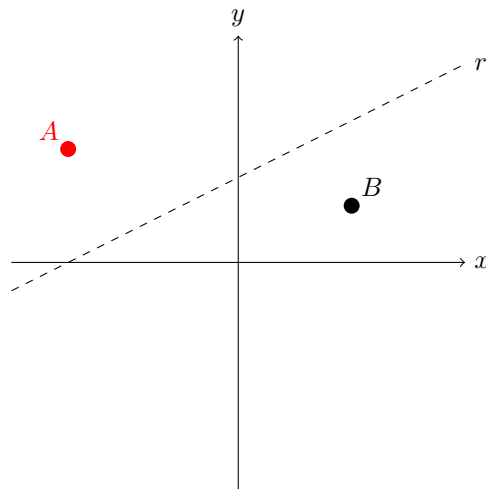
Il parametro della funzione non è altro che **l'equazione di una retta**.

In particolare, se consideriamo la retta che separa gli spazi del piano, vediamo che la retta è **capace di separare 2 punti nel piano**.

Cosa abbiamo:

- Un neurone
- 2 valori in input
- 2 pesi

Con la funzione di attivazione **sign** si avrà come output una retta che separa



dei punti.

Nota: Quando è utile avere un valore output che non è una separazione di punti in un piano? Nel caso della **regressione**.

Oltre la funzione sign:

- Funzione di attivazione lineare
- Funzione di attivazione sigmoid
- Funzione di attivazione Tanh

■ 3.4 Rappresentazione le funzioni logiche

:

◆ 3.4.1 AND

Immaginiamo di avere:

- $x_1, x_2 \in \{0, 1\}$
- Bias= -30
- Funzione di attivazione: Logistica o Sigmoid

Come facciamo a rappresentare un AND?

$$h(x) = g(-30 + 20x_1 + 20x_2) \quad (3)$$

x_1	x_2	$h(x)$
0	0	1
0	1	0
1	0	0
1	1	1

Lo stesso ragionamento vale per:

- OR
- NOT
- $(\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$

◆ 3.4.2 Il problema dello XOR

Il perceptrone non può imparare regioni che non sono linearmente separabili.

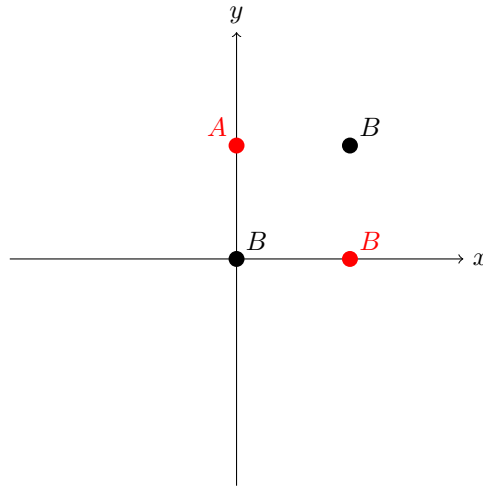


Figure 1: XOR non possibile con 1 perceptrone

Come vediamo qui non possiamo tracciare una retta per dividere i due punti. In questo caso ci serve una funzione **non lineare**.

La soluzione a questo problema è **aggiungere LAYER** alla rete neurale. Aggiungere un layer significa aggiungere un neurone successivamente ad un altro.

Maggiore è il numero di layer, maggiore diventa la potenza espressiva della rete. Praticamente possiamo catturare qualsiasi cosa aggiungendo layer alla rete. Questo è **il vero potere del deep learning**.

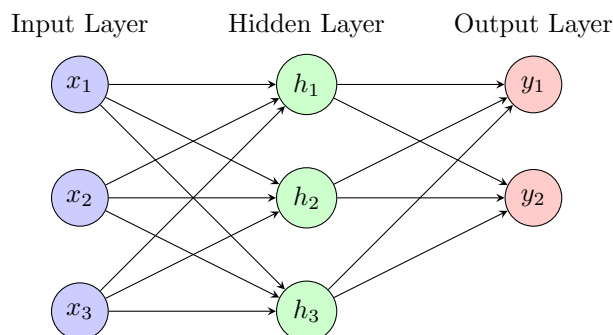


Figure 2: Neural Network con 1 hidden layer

■ 3.5 Gli ingredienti di una rete neurale

◆ 3.5.1 Il grafo g

Un grafo $g = \{N, E\}$ è un grafo diretto pesato con label.

Ogni nodo $i \in N$ viene chiamato **neurone** o **percetttrone**. Per ogni nodo ci sono 2 targhette:

- Un valore a_i che viene chiamato **attivazione**
- Una funzione di attivazione f_i che viene applicata all'attivazione, che produce un output z_i

Ogni arco $e = \{j \in N \rightarrow i \in N\} \in E$ associato con un peso $w_{j,i}$.

Ogni nodo i è anche coinvolto con un arco speciale, con un nodo fantasma, che viene chiamato **bias** b_i .

Nota: $z_i = f_i(a_i)$. E $a_i = b_i + \sum_{j:j \rightarrow i \in E} w_{j,i} z_j$

La combinazione dei neuroni connessi costruisce il grafo. I nodi che condividono gli stessi input sono raggruppati in **layers**.

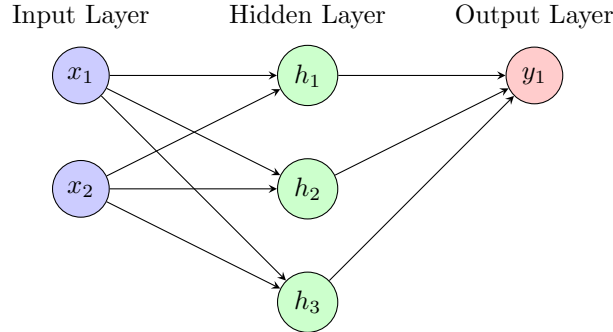


Figure 3: Neural Network con 1 hidden layer

Il risultato finale dipende da tutti i parametri del grafo, e anche dal tipo di funzione di attivazione che viene usata.

Ma come facciamo a scegliere il corretto valore dei parametri del grafo?

◆ 3.5.2 La funzione di loss

Il grafo è un operatore algebrico non lineare: $g(\vec{x}|W, B)$. In questo grafo non sappiamo il valore di **B** e **W**. La fase di apprendimenti di una rete consiste nel trovare il **migliore** valore per **W** e per **B**. Ma cosa intendiamo con **migliore**?

Formalmente, l'unico modo che abbiamo per definire la conezione di migliore, è quello di *approssimare al meglio la funzione che vogliamo trovare*.

Consideriamo il valore di una funzione F su x_i e il vero valore di y_i . Noi vogliamo minimizzare la differenza tra $F(x_i)$ e y_i .

$$\min_{W, B} \frac{1}{n} \sum_{i=1}^n \text{loss}(\vec{y}_i, g(x_i|W, B)) \quad (4)$$

con:

- $\text{loss}(\vec{y}_i, g(x_i|W, B))$ è una funzione che misura la differenza tra y_i e $g(x_i|W, B)$
- n è il numero di esempi
- W e B sono i parametri del grafo
- $g(x_i|W, B)$ è il valore di output del grafo
- \vec{y}_i è il valore di output vero

La funzione di **loss** dipende dal tipo di task che dobbiamo svolgere.

◆ 3.5.3 Funzione di loss per Regressione

Mean Absolute Error:

$$\frac{1}{n} \sum_{i=1}^n |y_i - g(\vec{x}_i|W, B)| \quad (5)$$

- Considera tutti gli errori con lo stesso peso
- Non è differenziabile in 0

Mean Squared Error:

$$\frac{1}{n} \sum_{i=1}^n (y_i - g(\vec{x}_i|W, B))^2 \quad (6)$$

- Gli errori più grandi hanno un peso maggiore
- È differenziabile in 0
- È più sensibile agli outliers

Domanda: quale si usa tra le due? Dall'approccio greedy, **usa entrambe**.
Esistono anche altre funzioni:

- Smooth Absolute Error
- Huber Loss

◆ 3.5.4 Funzione di loss per classificazione

Binary Cross Entropy [BCE] Viene usata se $y_i \in \{0, 1\}$, $g(\vec{x}_i|W, B) \in [0, 1]$.

$$BCE = -\frac{1}{n} \sum_{i=1}^n y_i \log(g(\vec{x}_i|W, B)) - (1 - y_i) \log(1 - g(\vec{x}_i|W, B)) \quad (7)$$

Categorical Cross Entropy [CCE]

$$CCE = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m y_{i,j} \log(g(\vec{x}_i|W, B)_j) \quad (8)$$

■ 3.6 L'ottimizzatore o

Come risolviamo il problema di:

$$\min_{W, B} \frac{1}{n} \sum_{i=1}^n \text{loss}(\vec{y}_i, g(\vec{x}_i|W, B)) \quad (9)$$

Il problema lo risolviamo calcolando il **gradiente** della funzione loss, lo poniamo uguale a 0 e controlliamo se siamo in un punto di **massimo**, **minimo** o **punto di sella**.

$$\frac{\partial loss}{\partial W} = 0 \quad (10)$$

Abbiamo bisogno di un metodo iterativo per trovare una soluzione. Ci vuole un'**euristica**. Tipicamente, siamo soddisfatti di un **minimo locale**, e si utilizza, appunto, il **metodo di discesa del gradiente**.

Differenza minimo locale e globale

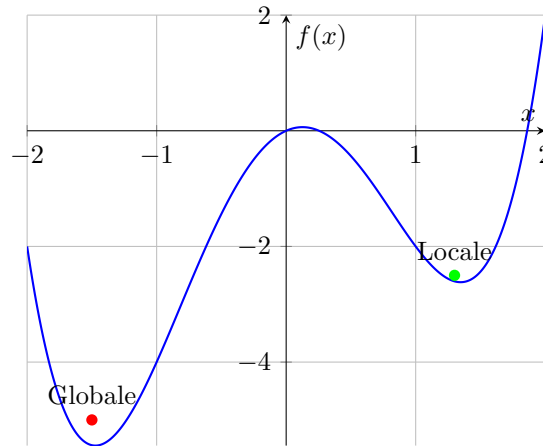


Figure 4: Differenza tra minimo locale e globale

◆ 3.6.1 Il metodo di discesa del gradiente

Sia $F(\vec{x})$ una funzione differenziabile.

- $F(\vec{x})$ decresce più veloce nella direzione del gradiente negativo
- $F(\vec{x})$ cresce più veloce nella direzione opposta al gradiente

$$a^{new} = a^{old} - \eta \cdot \nabla F(a^{old}) \quad (11)$$

Il parametro η viene chiamato **learning rate** e determina il comportamento dell'ottimizzazione. Da notare che e' l'unico **parametro**.

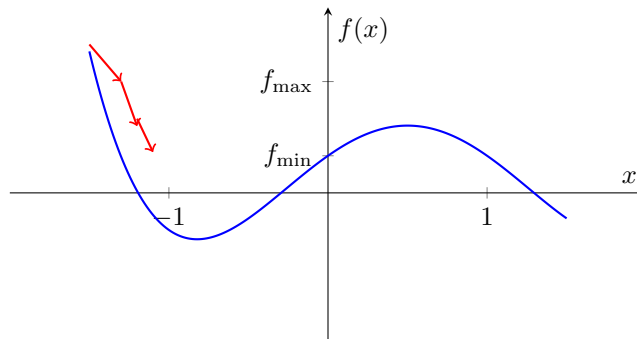


Figure 5: Discesa del gradiente

Ci sono ancora problemi: Questo metodo non è esente da problemi, quindi abbiamo:

- Dipendentemente dal punto di inizio, abbiamo un risultato diverso. Dobbiamo capire che, giustamente, se ci accontentiamo di un minimo locale, non avremo quasi mai lo stesso minimo per ogni allenamento della rete.

Nota: con **topologia** intendiamo il numero di layer e il numero di neuroni per layer.

Diciamo che in generale, gli steps per allenare una rete sono:

```

1   for each topology T in  $C_t$ 
2       for each  $\eta$  in  $C_\eta$ 
3           for each initialization I in  $C_I$ 
4               Applica la discesa del gradiente

```

◆ 3.6.2 Il metodo di discesa del gradiente stocastico

La differenza del metodo di discesa del gradiente stocastico è che, invece di calcolare il gradiente su tutti i dati, si calcola il gradiente su un sottoinsieme di dati.

Ciò che viene fatto, per ogni step, invece di prendere la derivata di ogni loss function e poi fare i calcoli, prendo **un sample del mio dataset**, tipicamente chiamato **batch**. Ogni volta che calcolo la discesa del gradiente, calcolo la derivata **NON PER TUTTO IL DATASET**, ma solamente del **batch**. Questo OVVIAMENTE non mi assicura che ottimizzare ogni batch mi ottimizza anche l'intero dataset, ma non c'è modo di lavorare sull'intero dataset, poiché questo è troppo grande e richiede troppo tempo.

Nonostante tutto, utilizzando questa tecnica **si ottengono risultati comunque accettabili**.

Il *workflow* allora cambia in:

```

1   for each topology T in  $C_t$ 
2       for each  $\eta$  in  $C_\eta$ 

```

```

3     for each initialization I in  $C_I$ 
4       for each batch size  $b \in C_b$ 
5         Applica la discesa del gradiente stocastica

```

◆ 3.6.3 Linee guida sul learning rate

Epoca: Un'epoca è un passaggio dell'approssimazione della funzione

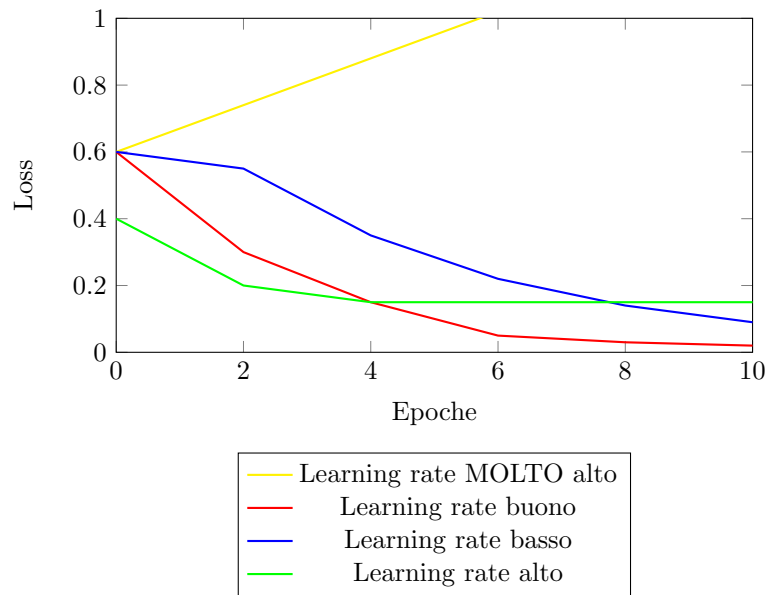


Figure 6: Esempio di andamento della funzione di loss in base al learning rate

Analisi dei learning rate:

- learning rate basso: troppo lento!
- learning rate alto: fermo con un minimo locale immediatamente
- learning rate MOLTO alto: non converge mai
- learning rate buono: decresce sempre e eventualmente converge

Nota: Il learning rate è un parametro che va **tunato**.

- Se il learning rate è troppo alto, non si riesce a convergere poiché si salta il minimo
- Se il learning rate è troppo basso, si rischia di non convergere mai

Il nostro obiettivo è quello di avere un learning rate che sia **giusto** e che permetta di avere un andamento della funzione di loss come quello in rosso, cioè avere una buona diminuzione della loss andando avanti con le epoche. Sempre decrescente e con la distanza minore dall'asse delle x.

Domanda: Vogliamo sempre la loss uguale a 0? **No.** Perché? Perché a differenza che in statistica dove vogliamo avere un errore il minimo. In deep learning non è così. Se la loss è 0, allora vuol dire che la rete ha imparato a memoria i dati, e non è quello che vogliamo. Vogliamo che la rete sia capace di generalizzare.

In particolare, noi **vogliamo fare predizioni**, soprattutto su **istanze ancora non viste**. Se avessimo una loss uguale a 0, probabilmente non saremmo in grado di avere delle predizioni decenti, poiché la rete non è in grado di generalizzare. Avendo la loss non a 0 rendiamo la rete capace di poter introdurre istanze ancora non viste in futuro.

Quindi, per statistica **loss = 0 : nice**, per **deep learning: insomma**.

■ 4 Classificazione Binaria

Descrizione di cosa andremo a fare: utilizziamo il dataset IMDb Dataset. Questo dataset è un Database compilato dagli utenti del sito. Le caratteristiche, in particolare, sono:

- 50.000 recensioni
- circa 50% positive e 50% negative
- 25.000 sono usate per il **training** e 25.000 sono usate per il **testing**. Anche queste sono il 50% positive e 50% negative.

La task che faremo su questo dataset prende il nome di **sentimental analysis**, ovvero analisi del sentimento.

Assunzione: Ciò che stiamo facendo ha senso solamente se assumiamo *che nel futuro avremo punti con una distribuzione abbastanza simile a quelli utilizzati per il training*.

La partizione di **test** è una partizione che non viene utilizzata per la fase di training, ma viene utilizzata successivamente per controllare se il modello si sta comportando bene nella predizione dei valori. Ci sono delle misure che hanno range $[0, 1]$ che ci permettono di capire quanto il modello si sta comportando bene.

Nota: il *test set* non viene fornito quando si lavora nel deep learning, altrimenti verrebbe usato per ottimizzare direttamente il modello. **Non si conosce inizialmente.**

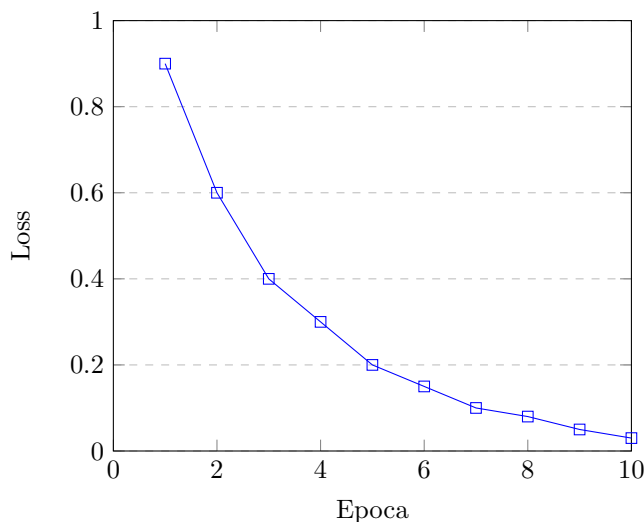


Figure 7: Grafico di una curva di loss che diminuisce all'aumentare delle epoche.

4.1 Caricamento e gestione del Dataset

```

1 from keras.datasets import imdb
2 (train_data, train_labels), (test_data, test_labels) = imdb.
   load_data(num_words=10000)

```

Listing 1: Caricamento del dataset IMDb Dataset.

- `training_data`: sono gli input del training
- `training_labels`: sono gli output del training
- `test_data`: sono gli input del test
- `test_labels`: sono gli output del test

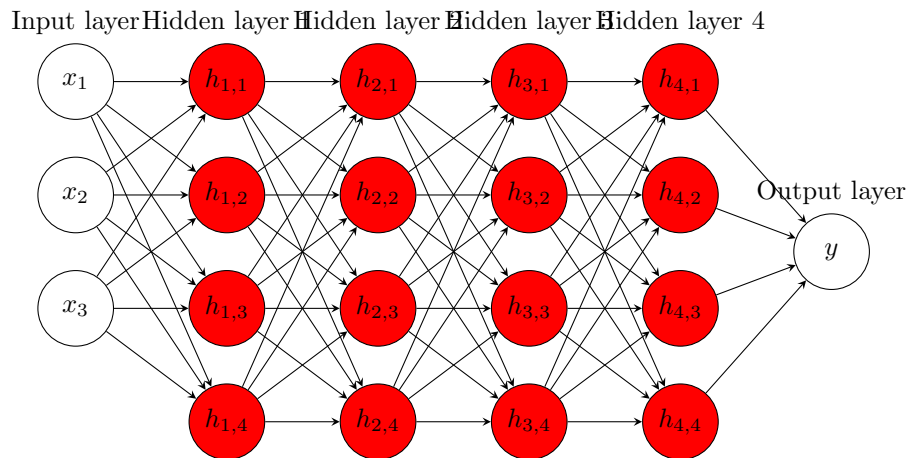


Figure 8: Neural network with 3 input nodes, 4 hidden layers with 3 nodes each, and 1 output layer with 1 node.

Quali sono i primi problemi che vanno gestiti in questo caso? Abbiamo il problema che i **dati sono parole e non numeri**. Quindi dobbiamo trasformare le parole in numeri.

Ogni parola viene **trasformato** in un numero. Ci troviamo però con delle recensioni che sono delle **sequenze di parole**, e quindi abbiamo delle **sequenze di numeri**. La soluzione è quella di usare un **set di parole codificate**. Mi spiego:

Immaginiamo di avere 10.000 parole encodeate. Salvo queste 10.000 parole in un array e associo ad ogni parola un indice di questo array.

0	1	2	3	4	5	6	7	8	9
---	---	---	---	---	---	---	---	---	---

Ora, cosa manca? **Manca l'ordine**. L'unica cosa che abbiamo è quindi una indicizzazione delle parole, ma non abbiamo alcuna informazione dell'ordine. **Manca anche totalmente la semantica.**

La struttura dell'input è quindi un **livello con 10.000 nodi di input**.
Input

```

1 import numpy as np
2 def vectorize_sequences(sequences, dimension=10000):
3     # Create an all-zero matrix of shape (len(sequences), dimension)
4     results = np.zeros((len(sequences), dimension))
5     for i, sequence in enumerate(sequences):
6         results[i, sequence] = 1. # set specific indices of results[i]
            to 1s
7     return results
8 # Our vectorized training data
9 x_train = vectorize_sequences(train_data)
10 # Our vectorized test data
11 x_test = vectorize_sequences(test_data)

```

Output

```

1 # Our vectorized labels
2 y_train = np.asarray(train_labels).astype('float32')
3 y_test = np.asarray(test_labels).astype('float32')

```

■ 4.2 Definizione della Rete Neurale

Codice iniziale della definizione della Rete

```

1 from keras import models
2 from keras import layers
3 model = models.Sequential()
4 model.add(layers.Dense(16, activation='relu', input_shape=(10000,)))
5 model.add(layers.Dense(16, activation='relu'))
6 model.add(layers.Dense(1, activation='sigmoid'))
7 model.compile(optimizer='rmsprop',
8               loss='binary_crossentropy',
9               metrics=['accuracy'])

```

Il **modello sequenziale** è quello migliore da dove iniziare. Cosa significa però modello sequenziale? Il modello sequenziale ha il concetto di **layer**. Ha la funzione **add** che aggiunge un layer sopra gli altri layer che sono già esistenti.

Primo Layer:

layers.Dense: un layer denso significa che ogni nodo di quel layer è collegato con ogni nodo del layer precedente. **16** è il numero di neuroni che si vogliono

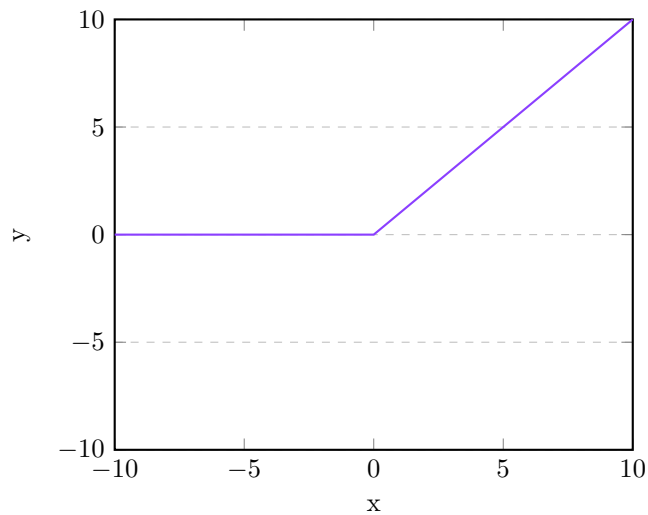
attivare in quel layer. `input_shape` è la dimensione dell'input. **10000** è la dimensione dell'input. Notare che si ha una **virgola** nell'input dopo il 10.000 e indica che *che stiamo aspettando una sequenza di vettori, ognuno di dimensione 10.000*. Cioè, praticamente stiamo dicendo **10.000** sono le parole per ogni recensione, e con la virgola stiamo dicendo che **non sappiamo quante recensioni abbiamo**. E' comodo quando non abbiamo un numero fisso di example del dataset da processare.

Quante connessioni ci sono?

$$16 \cdot 10000 + 16 = 160016$$

dove sono 16 i bias.

Cosa significa relu? ReLu è la funzione di attivazione.



E' una funzione di attivazione artificiale che fino a 0 è 0, e poi cresce linearmente. E' una funzione che si usa molto in deep learning.

Secondo layer:

Nel secondo layer abbiamo altri **16** nodi connessi con i precedenti 16 nodi. In tutto abbiamo

$$16 \cdot 16 + 16 = 272$$

connessioni, con 16 che sono i bias.

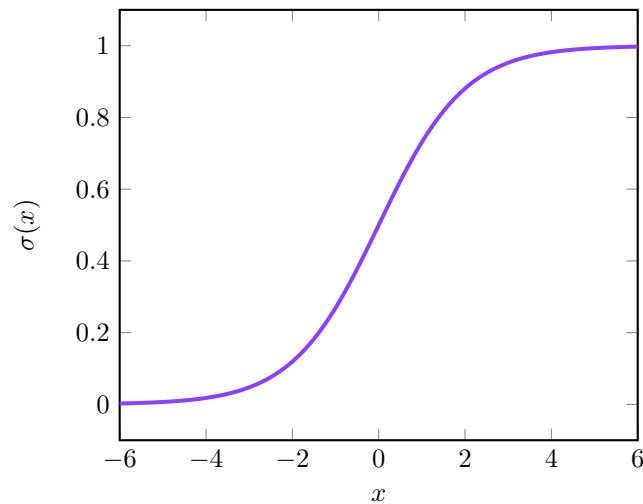
Layer di output:

In questo layer abbiamo un unico nodo connesso con i 16 nodi precedenti. In tutto abbiamo

$$16 \cdot 1 + 1 = 17$$

connessioni, con 1 che è il bias. La funzione di attivazione è la **sigmoid**, che è una funzione che ha range $[0, 1]$ che indica la probabilità che il dato appartenga alla classe 1.

Diciamo che è stato osservato che la *sigmoid è meglio nell'output*. Perché? Eh funziona così a quanto pare lol.



Compilazione del modello

In questa sezione si definisce il **loss** e l'**optimizer**.

L'ottimizzatore sarebbe, praticamente, la *discesa del gradiente*. Ci sono vari ottimizzatori ed è molto buono per te stesso provare gli ottimizzatori e vedere quale funziona meglio per il tuo problema. Letteralmente il deep learning :) Per quanto riguarda la loss abbiamo visto nella scorsa sezione le funzioni di loss che conosciamo.

La **metrica** è quella che viene usata per valutare effettivamente il modello. In questo caso si usa l'**accuracy** che praticamente è la percentuale di classificazioni corrette.

■ 4.3 Plotting del modello

E' un plotting molto base e non molto fancy, ma fa il suo lavoro diciamo

```
1 from keras.utils import plot_model
2 plot_model(model, show_shapes=True, show_layer_names=True)
```

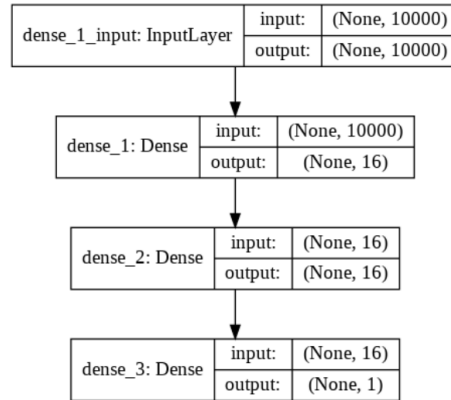


Figure 9: Plotting del modello.

■ 4.4 Training del modello e valutazione

Entra in gioco il concetto di **validation set**.

◆ 4.4.1 Validation Set

Quando ho un modello, come faccio a sapere se la rete che sto allenando si sta allenando in modo corretto? Tutto questo sapendo che **non abbiamo accesso al test set**. Pensiamo al **training set** e pensiamo ad un altro **split** interno al training set:

- Training set parziale
- Validation set

Quindi se avessimo 50.000 di grandezza del dataset:

- 25.000 test set
- 25.training set
 - 10.000 validation set
 - 15.000 training set parziale

```

1 x_val = x_train[:10000]
2 partial_x_train = x_train[10000:]
3 y_val = y_train[:10000]
4 partial_y_train = y_train[10000:]

```

Praticamente è come se usassimo il validation come se fosse un test set. Quindi andremo a plottare 2 curve:

- La loss sulle epoche
- L'accuracy sulle epoche

La loss ci da un'idea di quanto il modello si sta allenando bene. Se ci sono problemi, la figura è strana e non segue un andamento corretto, si modifica. Ma la cosa importante è la **validation accuracy**, che DEVE essere crescere in modo monotono e deve avvicinarsi a 1 il più possibile.

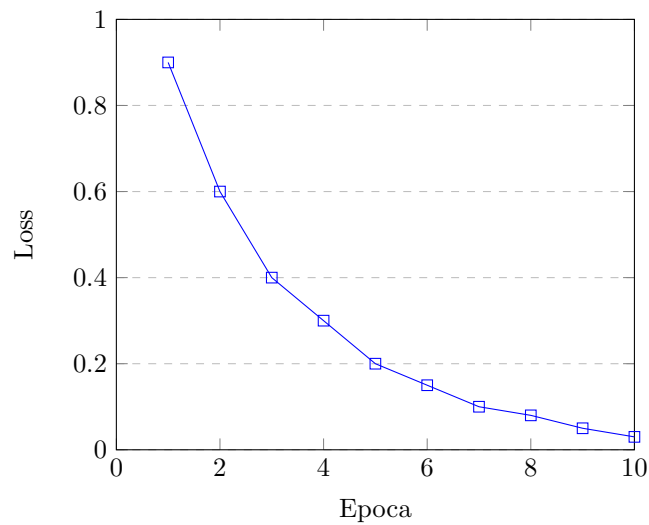


Figure 10: Grafico di una curva di loss che diminuisce all'aumentare delle epoche.

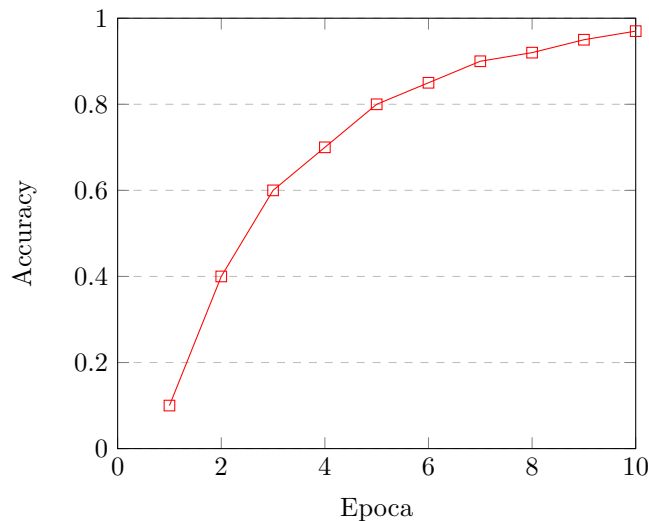


Figure 11: Grafico di una curva di accuracy che aumenta all'aumentare delle epoche.

Ma anche se avessimo queste curve in questa conformazione, **ancora non ci dice niente**. Il validation è in un senso insensato per l'output. Ciò che importerà sarà il **test set** che non è stato usato per il training.

Training code

```
1 history = model.fit(partial_x_train,
2                     partial_y_train,
3                     epochs=20,
4                     batch_size=512,
5                     validation_data=(x_val, y_val))
```

Plotting delle curve

Notare che questo è un plotting molto base e non molto fancy, e in futuro utilizzeremo **Tensorboard** che aggiornerà in tempo reale i grafici

```
1 import matplotlib.pyplot as plt
2 loss = history.history['loss']
3 val_loss = history.history['val_loss']
4 epochs = range(1, len(loss) + 1)
5 # "bo" is for "blue dot"
6 plt.plot(epochs, loss, 'bo', label='Training loss')
7 # b is for "solid blue line"
8 plt.plot(epochs, val_loss, 'b', label='Validation loss')
9 plt.title('Training and validation loss')
10 plt.xlabel('Epochs')
11 plt.ylabel('Loss')
12 plt.legend()
13 plt.show()
```

Notare che *history* è un dizionario e si può accedere come indicato a riga 1 e 2 come accedere a *loss* e *val_loss*.

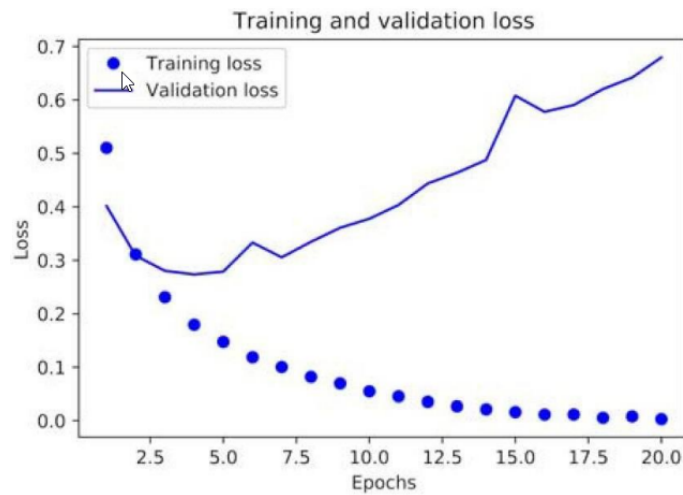


Figure 12: Rete sbagliata

Un grafico del genere ci mostra una rete che **funziona male!**.

```

1 plt.clf() # clear figure
2 acc = history_dict['binary_accuracy']
3 val_acc = history_dict['val_binary_accuracy']
4 plt.plot(epochs, acc, 'bo', label='Training acc')
5 plt.plot(epochs, val_acc, 'b', label='Validation acc')
6 plt.title('Training and validation accuracy')
7 plt.xlabel('Epochs')
8 plt.ylabel('Accuracy')
9 plt.legend()
10 plt.show()

```

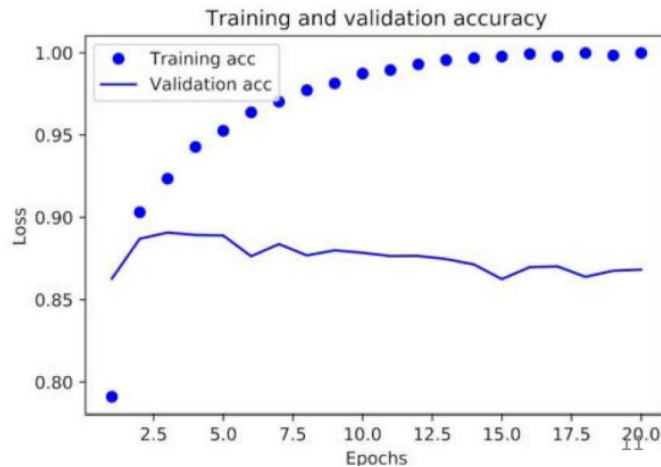


Figure 13: Overfitting

In questo caso la rete ha un problema di **overfitting!**.

4.5 Prediction

```
1 model.predict(x_test)
array([[0.91966152], [0.86563045], [0.99936908], ..., [0.45731062], [0.0038014], [0.79525089]], dtype = float32)
```

4.6 Come risolvere problemi di accuracy bassa?

Se quando andiamo a fare la validazione della nostra rete, come sappiamo come modificare qualcosa per rendere la rete migliore? Ci sono alcuni approcci per farlo:

- Cambiare la topologia
- Diminuire le epoche
- Cambiare il numero di nodi per qualche layer
- FORSE cambiare anche il learning rate

Si noti che facendo questo vanno interpretati i grafici per capire come sta andando la rete. Quando cominciamo a vedere che la **validation accuracy NON STA SOLAMENTE CRESCENDO** ma ci sono dei punti in cui diminuisce, siamo sicuri che **c'è qualche problema di mezzo**. Bisogna capire anche su cosa lavorare. Se cambiando i dati spesso la **validation accuracy** mostra problemi, allora bisogna lavorare per quello e cercare di trovare un modo per fare in modo che il problema non si presenti.

RECAPPONE: Se notiamo che nella LOSS ci sono problemi, possiamo tunare il **learning rate** e altri parametri per cercare di risolvere questi problemi. Ma la cosa più importante è la **metrica**. Se notiamo che la **metrica non rispecchia un andamento di effettivo apprendimento**, capiamo che la rete non si sta comportando nel modo corretto e non sta funzionando.

■ 4.7 Early Stopping

E' un metodo che consiste nel vedere quando **la loss** comincia ad avere dei comportamenti sospetti. Quando si **ferma** la rete in un dato momento, si impedisce alla rete di **overfittare** nella maggior parte dei casi.

■ 5 Classificazione Multiclasse

Questa sezione sarà molto corta, poiché praticamente è la stessa cosa della classificazione binaria, ma il layer finale della rete non sarà composto da un solo nodo, ma da **un numero di nodi pari al numero di classi da classificare**.

■ 5.1 Descrizione del dataset - Reuters

Il dataset che verrà utilizzato è il **Reuters Dataset**, che è un dataset di **news** che sono state classificate in **46 categorie**. Ogni categoria ha almeno 10 esempi nel training set.

Preprocessing dell'input

```

1 import numpy as np
2 def vectorize_sequences(sequences, dimension=10000):
3     results = np.zeros((len(sequences), dimension))
4     for i, sequence in enumerate(sequences):
5         results[i, sequence] = 1.
6     return results
7 # Our vectorized training data
8 x_train = vectorize_sequences(train_data)
9 # Our vectorized test data
10 x_test = vectorize_sequences(test_data)
```

◆ 5.1.1 Come processiamo l'output?

Nel deep learning, **gli output categorici** non vengono mai mappati ad una scala numerica. Si utilizza una tecnica che si chiama **one hot encoding**.

One Hot Encoding:

Consiste nel creare tante variabili numeriche quanti sono i valori della variabile categorica, e per quel dato example viene assegnato

- 1: se il valore della variabile categorica è quello

- 0: se il valore della variabile categorica non è quello, quindi in tutti gli altri casi

Da notare che il one hot encoding viene fatto sia per le labels di output, ma non è sbagliato farlo anche per gli attributi in alcuni casi.

```

1 def to_one_hot(labels, dimension=46):
2     results = np.zeros((len(labels), dimension))
3     for i, label in enumerate(labels):
4         results[i, label] = 1.
5     return results
6 # Our vectorized training labels
7 one_hot_train_labels = to_one_hot(train_labels)
8 # Our vectorized test labels
9 one_hot_test_labels = to_one_hot(test_labels)
10
11 #OPPURE ALLO STESSO MODO
12
13 from keras.utils.np_utils import to_categorical
14
15 one_hot_train_labels = to_categorical(train_labels)
16 one_hot_test_labels = to_categorical(test_labels)

```

Nota: nel one hot encoding DOBBIAMO AVERE 1 SOLO VALORE per il dato example, e tutti gli altri non devono essere attivi. La domanda è, quindi, **quale activation function usiamo?**

Immaginiamo questa situazione: Abbiamo y_1, y_2, y_3 che sono i valori di output. Vogliamo solamente uno dei tre. Se normalizzassimo i dati in questo modo:

- $p_1 = \frac{y_1}{y_1 + y_2 + y_3}$
- $p_2 = \frac{y_2}{y_1 + y_2 + y_3}$
- $p_3 = \frac{y_3}{y_1 + y_2 + y_3}$

E se sommiamo $p_1 + p_2 + p_3$ otteniamo 1. Quindi, in questo caso, possiamo usare la **softmax** come activation function.

◆ 5.1.2 Softmax activation function

Se prendiamo un vettore $\vec{y} = (y_1, y_2, y_3)$, la softmax è definita come:

$$S(y_i) = \frac{e^{y_i}}{\sum_{j=1}^3 e^{y_j}} \quad (12)$$

In output abbiamo: $\vec{p} = (p_1, p_2, p_3)$, dove p_i è la probabilità che il dato example appartenga alla classe i .

Nota: quando dividiamo qualcosa, abbiamo **sempre** un problema riguardo la stabilità numerica. La divisione è molto critica, poiché **potrebbe essere vicina a 0**. Sappiamo che se usiamo $e_1^y + e_2^y + e_3^y$ difficilmente si avvicina a 0.

```
1 from keras import models
2 from keras import layers
3 model = models.Sequential()
4 model.add(layers.Dense(64, activation='relu', input_shape
5   =(10000,)))
6 model.add(layers.Dense(64, activation='relu'))
7 model.add(layers.Dense(46, activation='softmax'))
8 model.compile(optimizer='rmsprop',
9               loss='categorical_crossentropy',
10              metrics=['accuracy'])
```

Nota: Dobbiamo avere un numero di nodi di output quanto il numero di classi, ma **altra cosa**, il numero di nodi interni deve essere sicuramente maggiore di 46, altrimenti ci sarebbe una perdita di informazioni.

Nota 2: La funzione di attivazione dell'ultimo layer è una **softmax**.

Validation set e Training set

```
1
2 x_val = x_train[:1000]
3 partial_x_train = x_train[1000:]
4 y_val = one_hot_train_labels[:1000]
5 partial_y_train = one_hot_train_labels[1000:]
6 history = model.fit(partial_x_train,
7                    partial_y_train,
8                    epochs=20,
9                    batch_size=512,
10                   validation_data=(x_val, y_val))
```

Loss

```
1 import matplotlib.pyplot as plt
2 loss = history.history['loss']
3 val_loss = history.history['val_loss']
4 epochs = range(1, len(loss) + 1)
5 plt.plot(epochs, loss, 'bo', label='Training loss')
6 plt.plot(epochs, val_loss, 'b', label='Validation loss')
7 plt.title('Training and validation loss')
8 plt.xlabel('Epochs')
9 plt.ylabel('Loss')
10 plt.legend()
11 plt.show()
```

Accuracy

```
1 plt.clf() # clear figure
2 acc = history.history['acc']
3 val_acc = history.history['val_acc']
4 plt.plot(epochs, acc, 'bo', label='Training acc')
5 plt.plot(epochs, val_acc, 'b', label='Validation acc')
6 plt.title('Training and validation accuracy')
7 plt.xlabel('Epochs')
8 plt.ylabel('Acc')
9 plt.legend()
10 plt.show()
```

Nota: L'accuracy è importante dipendentemente dall'utilizzo che bisogna farne. Ad esempio, *in ambito medico* vogliamo **minimizzare i falsi negativi**. Questo implica che la misura che usiamo dipende dall'applicazione che se ne fa.

Capitoli di laboratorio

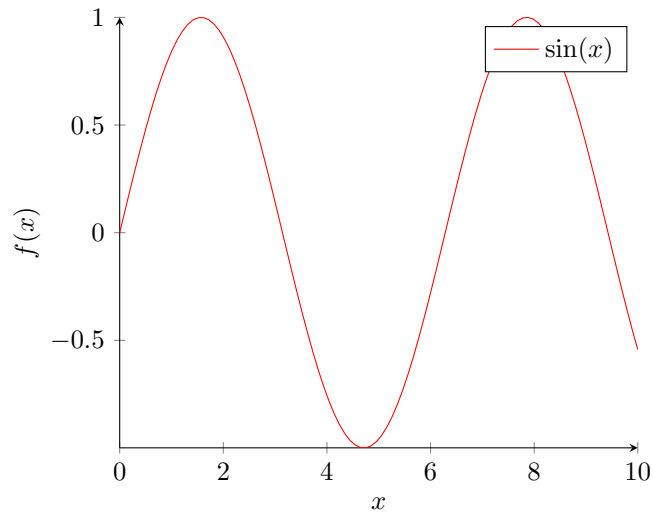
■ 6 Lab: Introduzione Python

■ 6.1 Matplotlib

◆ 6.1.1 Plots

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 x = np.linspace(0, 10, 100)
5 plt.plot(x, np.sin(x))
6 plt.show()
```

Non c'è molto da dire, il codice è autoesplicativo. La funzione `plot` prende in input due array, uno per l'asse delle ascisse e uno per l'asse delle ordinate. In questo caso, `x` è un array di 100 punti equidistanti tra 0 e 10, mentre `np.sin(x)` è un array di 100 punti che rappresentano il seno dei punti di `x`. La funzione `show` mostra il grafico.

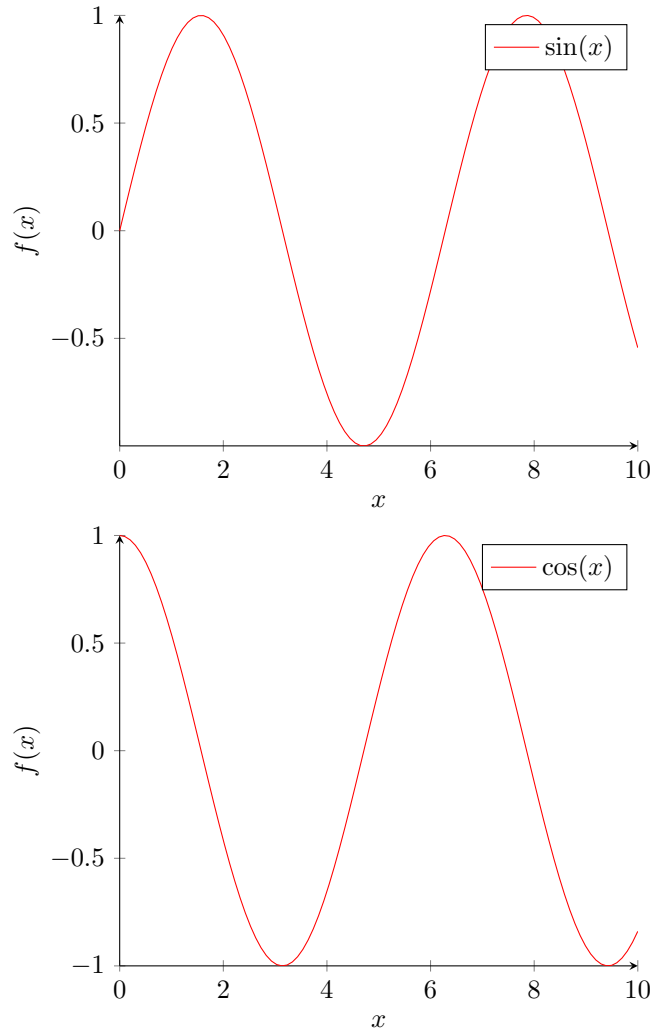


◆ 6.1.2 Sub-plots

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 x = np.linspace(0, 10, 100)
5 plt.subplot(2, 1, 1)
6 plt.plot(x, np.sin(x))
```

```
7 plt.subplot(2, 1, 2)
8 plt.plot(x, np.cos(x))
9 plt.show()
```

La funzione `subplot` prende in input tre parametri: il numero di righe, il numero di colonne e l'indice del subplot corrente. Nel caso di questo esempio, il subplot corrente è il primo, quindi viene mostrato il grafico del seno. Poi viene mostrato il secondo subplot, che è quello del coseno.



C'è anche qui poco da dire, il codice è autoesplicativo.

■ 6.2 NumPy

La libreria NumPy è una libreria per Python che permette di lavorare con array multidimensionali. Per importare la libreria, basta scrivere `import numpy as`

`np.`

Mi rompo le palle in maniera assurda di scrivere tutti gli esempi. Quindi questo capitolo penso sia abbastanza inutile.

E' possibile trovare il codice di **numPy** a questo link: www.ciao.it

■ 7 Lab: Reti Neurali da zero

■ 7.1 Introduzione

Partiamo dicendo una cosa molto importante: **per quale motivo usiamo la backpropagation?**

$$\frac{y - b}{x} = w \quad (13)$$

Ma possiamo scriverlo come:

$$(y - b) \cdot x^{-1} = w \quad (14)$$

Ora però, se consideriamo:

- y il vettore risultante
- w la matrice dei pesi
- x il vettore di input

Calcolare l'inversa di \mathbf{x} non è una cosa così poco costosa, anzi. Per questo motivo utilizziamo il training delle reti come la backpropagation e la discesa del gradiente.

Ora, andiamo più a fondo. Facciamo un esempio più pratico.

■ 7.2 Esempio pratico

Immaginiamo di avere un dataset di 2 features e una label da identificare.

x_1	x_2	y
1	2	0
2	3	0
3	4	0
4	5	0
5	6	0
6	7	1
7	8	1
8	9	1
9	10	1
10	11	1

Table 1: Dataset di esempio

Ogni livello di una rete neurale può essere rappresentato attraverso le **matrici**.

Ma passiamo alla definizione formale:

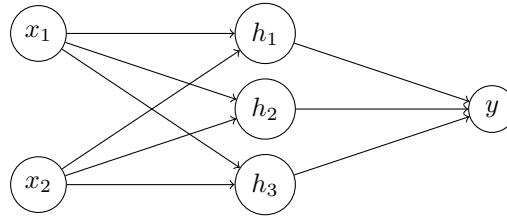


Figure 14: Rete neurale di esempio

Definizione Rete Neurale: Una rete neurale è una tupla:

$$NN = \{g, l, o, i, fpp\} \quad (15)$$

con:

- g : il grafo
- l : la funzione loss
- o : l'ottimizzatore
- i : l'inizializzatore
- fpp : la fix point procedure

Nota 1: l'ottimizzatore intende in quale modo si performa la **discesa del gradiente**.

Nota 2: l'inizializzatore è la funzione che inizializza i pesi della rete neurale. Ci possono essere diversi algoritmi per gli inizializzatori, ma non c'è modo di sapere quale funziona meglio, poiché non si può sapere nelle reti neurali.

◆ 7.2.1 Il grafo

Il solito grafo che abbiamo visto in precedenza è un grafo aciclico diretto. Ha dei nodi che sono i percettroni, che hanno degli archi composti dai pesi, che hanno una funzione di attivazione, che prendono in input un valore e ne sputano fuori uno chiamando la funzione di attivazione sull'input.

Nota: la funzione di attivazione è una funzione non lineare.

◆ 7.2.2 La funzione loss

L'obiettivo della funzione loss è quello di misurare la distanza tra il valore predetto e il valore reale.

$$L(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2 \quad (16)$$

◆ 7.2.3 L'ottimizzatore

L'ottimizzatore ci permette di trovare una soluzione ottimale **non ottima**, cioè trovare **il minimo della funzione di loss**. La tecnica che si usa è quella della **discesa del gradiente**.

◆ 7.2.4 Discesa del gradiente

La discesa del gradiente sfrutta un parametro chiamato **learning rate** che permette di capire quanto la discesa del gradiente deve essere veloce. Se il learning rate è troppo alto, la discesa del gradiente potrebbe non convergere, se è troppo basso, la discesa del gradiente potrebbe convergere troppo lentamente.

Nota: il learning rate è un ottimizzatore *molto naive*.

◆ 7.2.5 Inizializzatore

Il metodo di inizializzazione può cambiare di molto il risultato della rete neurale. In pratica assegna un valore ai *pesi* e ai *bias*. Alcuni metodi sono:

- Inizializzazione a 0: Questo ha alcuni problemi, perché non si ha diversificazione tra i nodi e i nodi nascosti diventano simmetrici.
- inizializzazione costante: stessi problemi della precedente
- Inizializzazione Random: Si fa seguendo una distribuzione uniforme oppure una distribuzione normale.

◆ 7.2.6 Fix Point Procedure

La procedura per il training di una rete neurale è un processo che si basa su alcuni steps:

```

1
2 net CustomNeuralNetwork(...)
3 initialize_weights_and_biases (net)
4 optimizer = myOptimizer(...)
5 loss_function = myLoss Function(...)
6 epochs = ... # the number of dataset scans
7 history = [] # a list containing the loss evolution
8 for epoch in range(epochs):
9     optimizer.reset() # it may have an internal status
10    loss = loss_fuction (out_target, net(input)) back_propagation (
        loss, optimizer, net)
11    history.append(loss)

```

■ 7.3 Esempio da zero

◆ 7.3.1 La funzione Sigmoid

Spendiamo qualche parola sulla funzione sigmoid, che è una funzione molto importante per le reti neurali.

In particolare, la funzione avrà un valore di attivazione compreso tra 0 e 1 e, in particolare, quando l'input è 0, la funzione ha valore 0.5. Se ha un valore basso, sarà zero e chiaramente se sarà alto avrà valore 1.

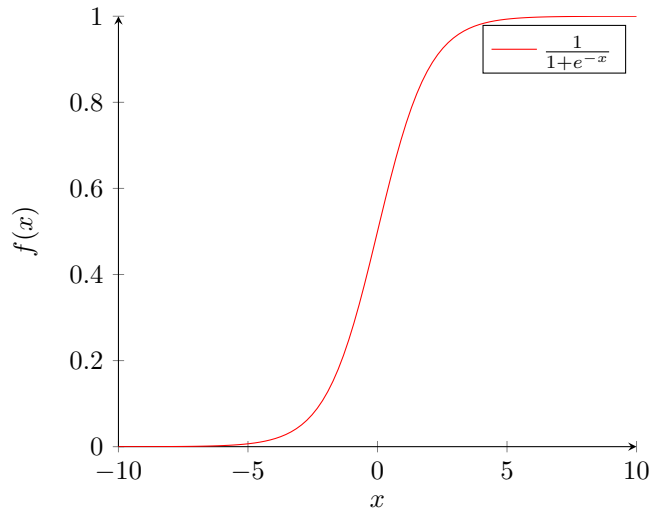


Figure 15: Funzione Sigmoid

Detto questo, **per quale motivo è importante?** Se pensiamo all'*inizializzazione*, che tipi di valori è meglio avere? Avere tutti i valori a zero porterebbe a valori simmetrici e quindi a problemi di convergenza. Vogliamo dei valori di pesi che siano **distribuiti uniformemente intorno allo 0**.

■ 7.4 Tangente Iperbolica TanH

La funzione TanH è una funzione che ha un valore di attivazione compreso tra -1 e 1 e, in particolare, quando l'input è 0, la funzione ha valore 0. Se ha un valore basso, sarà -1 e chiaramente se sarà alto avrà valore 1.

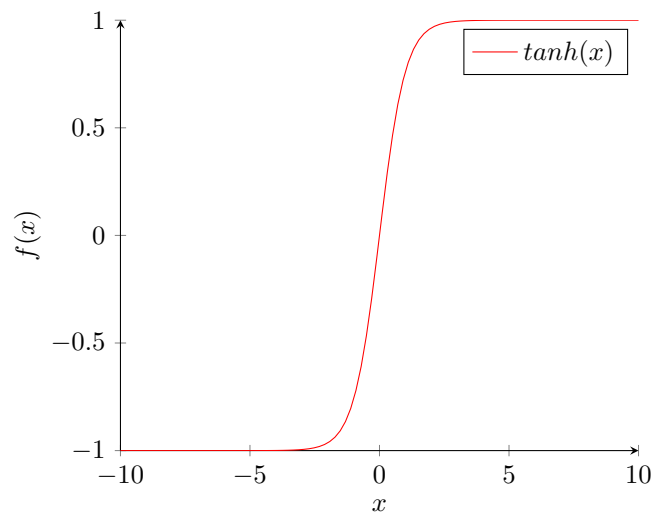


Figure 16: Funzione TanH

■ 7.5 ReLu

La funzione ReLu è una funzione che ha un valore di attivazione pari a 0 quando l'input è negativo e ha un valore di attivazione pari all'input quando l'input è positivo.

Parole di Adornetto: E' buona? Si. E' stabile? Si. Perché? Boh.

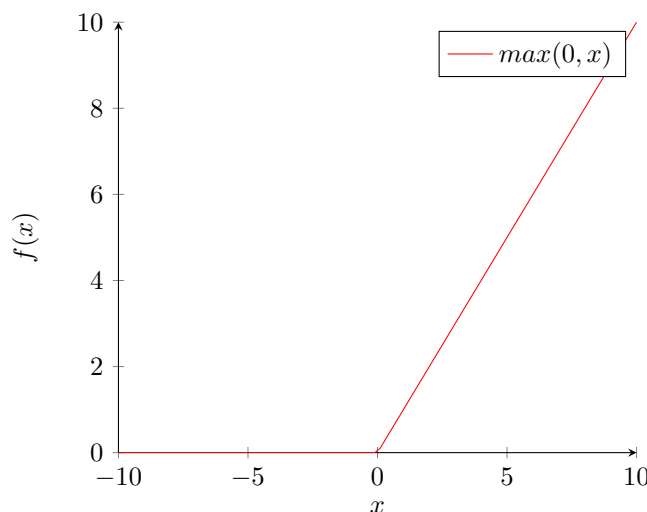


Figure 17: Funzione ReLu

7.6 Scalare i valori

Skip molto avanti riguardo l'esempio, ma si parla di scalare i dati.

Lo scaling dei dati è un'operazione importante nel machine learning perché i dati possono essere su scale diverse e questo può causare problemi durante l'allenamento del modello. Ad esempio, se abbiamo due variabili di input, una che varia da 0 a 1 e l'altra che varia da 0 a 1000, la seconda variabile avrà un impatto molto maggiore sull'output del modello rispetto alla prima. Ciò può portare a problemi come l'overfitting, in cui il modello si adatta troppo ai dati di addestramento e non generalizza bene sui dati di test.

Per risolvere questo problema, si utilizzano tecniche di scaling dei dati per portare tutte le variabili su una scala comune. Ci sono diverse tecniche di scaling, come la normalizzazione e la standardizzazione, che possono essere utilizzate a seconda del tipo di dati e del modello utilizzato.

Le funzioni di attivazione, come la funzione ReLU, possono anche essere influenzate dalla scala dei dati di input. Se i dati non sono scalati correttamente, la funzione di attivazione potrebbe produrre valori che non permettono un allenamento corretto del modello. Lo scopo dello scaling dei dati è quello di mantenere i dati intorno allo 0, in modo che le funzioni di attivazione possano produrre valori che permettono un allenamento corretto del modello.

- Standardizzazione: La standardizzazione è una tecnica di scaling dei dati che assume che i dati siano distribuiti normalmente all'interno di ogni feature e li scala in modo che la distribuzione abbia una media uguale a 0 e una deviazione standard uguale a 1. Questa tecnica funziona bene

quando i dati hanno una distribuzione normale, ma non funziona bene se i dati hanno una distribuzione non normale.

- Normalizzazione: La normalizzazione è una tecnica di scaling dei dati che scala i valori di ogni feature in modo che siano compresi tra 0 e 1. Questa tecnica funziona bene quando i valori di input non hanno una distribuzione normale. Ad esempio, se i valori di input sono compresi tra 0 e 1000, la normalizzazione li porterà su una scala compresa tra 0 e 1.

■ 7.7 Plottare la loss

La funzione loss è una misura dell'errore del modello durante l'allenamento. L'obiettivo dell'allenamento è quello di minimizzare la funzione loss, ovvero di ridurre l'errore del modello. Durante l'allenamento, il modello viene eseguito su un set di dati di addestramento e la funzione loss viene calcolata per ogni esempio di addestramento. L'errore totale del modello è la somma di tutte le funzioni loss per ogni esempio di addestramento.

L'allenamento del modello avviene in epoche, ovvero in cicli di esecuzione su tutti i dati di addestramento. L'obiettivo è che la funzione loss diminuisca ad ogni epoca, ovvero che l'errore del modello diminuisca man mano che il modello viene addestrato su più dati. Se la funzione loss non diminuisce ad ogni epoca, significa che il modello non sta imparando abbastanza dai dati di addestramento e potrebbe essere necessario modificare l'architettura del modello o i parametri di addestramento.

Se la funzione di loss arriva a 0, vuol dire che **siamo in minimo globale**, che non è comune.

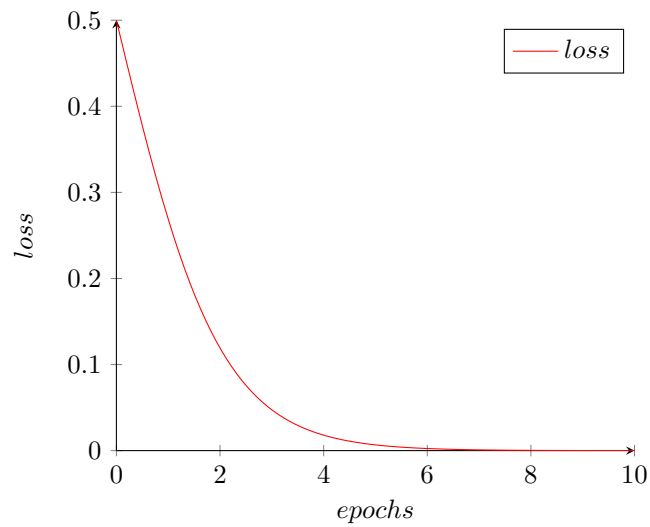


Figure 18: Funzione loss

Se invece finiamo con un grafico che arriva fino a 0.1 e poi si stabilizza, potremmo essere in un **minimo locale**, oppure semplicemente la topologia della Rete permette di avere questo risultato come massimo risultato.