

UNIVERSITÀ DI BOLOGNA

Attività progettuale di Sistemi Intelligenti M

Reti neurali artificiali: dal MLP agli ultimi modelli di CNN

Autore:

Ali Alessio Salman

22 agosto 2017

Indice

1	Reti Neurali Artificiali: le basi	1
1.1	Breve introduzione	1
1.2	Multi-layer Perceptron	2
1.2.1	Strati Nascosti	3
1.3	Caso di studio: prevedere il profitto di un ristorante	3
1.4	Implementazione da zero di un MLP	4
1.4.1	Dataset e Architettura	4
1.4.2	Forward Propagation	6
1.4.3	Backpropagation	7
1.4.4	Verifica numerica del gradiente	8
1.4.5	Addestramento	8
1.5	Ottimizzazione: diverse tecniche	8
1.6	Overfitting: come rilevarlo e risolverlo	8
1.7	Risultati	8
2	Reti Neurali Convolutionali	9
2.1	Breve introduzione	9
2.2	Architettura	9
2.2.1	Strato di Convoluzione	9
2.2.2	Strato di Pooling	9
2.2.3	Strato completamente connesso (FC)	9
2.3	Applicazioni e risultati	9
3	CNN: implementazione e addestramento	11
3.1	Il framework Torch	11
3.2	Modello di CNN basato su LeNet	11
3.2.1	Strato di Convoluzione	11
3.2.2	Strato di Pooling	11
3.2.3	Strato completamente connesso (FC)	11
3.3	Dataset d'esempio: CIFAR	11
3.4	Addestramento su CIFAR	11
4	Addestrare un modello allo stato dell'arte	13
4.1	ResNet	13
4.2	Implementazione di Resnet di Facebook	13
4.3	Addestramento su CIFAR	13
4.4	LeNet vs Resnet: confronto	13
5	Caso d'uso attuale: fine-tuning su dataset arbitrario	15
5.1	Dataset	15
5.2	Fine-Tuning	15
5.3	ModelZoo	15
5.3.1	Fine-tuning su Resnet	15

6 Conclusioni	17
A Appendice A: Codice aggiionale ed il framework Torch	19

Capitolo 1

Reti Neurali Artificiali: le basi

1.1 Breve introduzione

Una rete neurale artificiale – chiamata normalmente solo rete neurale (in inglese *Neural Network*) – è un modello di calcolo adattivo, ispirato ai principi di funzionamento del sistema nervoso degli organismi evoluti che secondo l’approccio connessionista [WConnessionismo] possiede una complessità non descrivibile con i metodi simbolici. La caratteristica fondamentale di una rete neurale è che essa è capace di acquisire conoscenza modificando la propria struttura in base alle informazioni esterne (i dati in ingresso) e interne (tramite le connessioni) durante il processo di apprendimento. Le informazioni sono immagazzinate nei parametri della rete, in particolare, nei pesi associati alle connessioni. Sono strutture non lineari in grado di simulare relazioni complesse tra ingressi e uscite che altre funzioni analitiche non sarebbero in grado di fare.

L’unità base di questa rete è il neurone artificiale introdotto per la prima volta da McCulloch e Pitts nel 1943 (fig. 1.1).

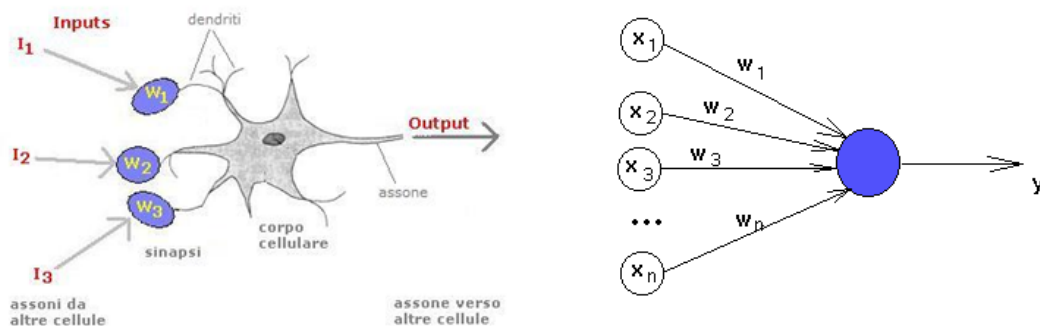


Figura 1.1: Modello di calcolo di un neurone (a sinistra) e schema del neurone artificiale (a destra)

Si tratta di un’unità di calcolo a N ingressi e 1 uscita. Come si può vedere dall’immagine a sinistra gli ingressi rappresentano le terminazioni sinaptiche, quindi sono le uscite di altrettanti neuroni artificiali. A ogni ingresso corrisponde un peso sinaptico w , che stabilisce quanto quel collegamento sinaptico influisca sull’uscita del neurone. Si determina quindi il potenziale del neurone facendo una somma degli ingressi, pesata secondo i pesi w .

A questa viene applicata una funzione di trasferimento non lineare:

$$f(x) = H\left(\sum_i (w_i x_i)\right) \quad (1.1)$$

ove H è la funzione gradino di Heaviside [**WHeaviside**]. Vi sono, come vedremo, diverse altre funzioni non lineari tipicamente utilizzate come funzioni di attivazioni dei neuroni. Nel '58 Rosenblatt propone il modello di *Percettrone* rifinendo il modello di neurone a soglia, aggiungendo un termine di *bias* e un algoritmo di apprendimento basato sulla minimizzazione dell'errore, cosiddetto *error back-propagation* [**WPercettrone**].

$$f(x) = H\left(\sum_i (w_i x_i) + b\right), \quad \text{ove } b = \text{bias} \quad (1.2)$$

$$w_i(t+1) = w_i(t) + \eta \delta x_i(t) \quad (1.3)$$

dove η è una costante di apprendimento strettamente positiva che regola la velocità di apprendimento, detta *learning rate* e δ è la discrepanza tra l'output desiderato e l'effettivo output della rete.

Il percettrone però era in grado di imparare solo funzioni linearmente separabili. Una maniera per oltrepassare questo limite è di combinare insieme le risposte di più percettroni, secondo architetture multistrato.

1.2 Multi-layer Perceptron

Il Multi-layer Perceptron (*MLP*) o percettrone multi-strato è un tipo di rete feed-forward che mappa un set di input ad un set di output. È la naturale estensione del percettrone singolo e permette di distinguere dati non linearmente separabili.

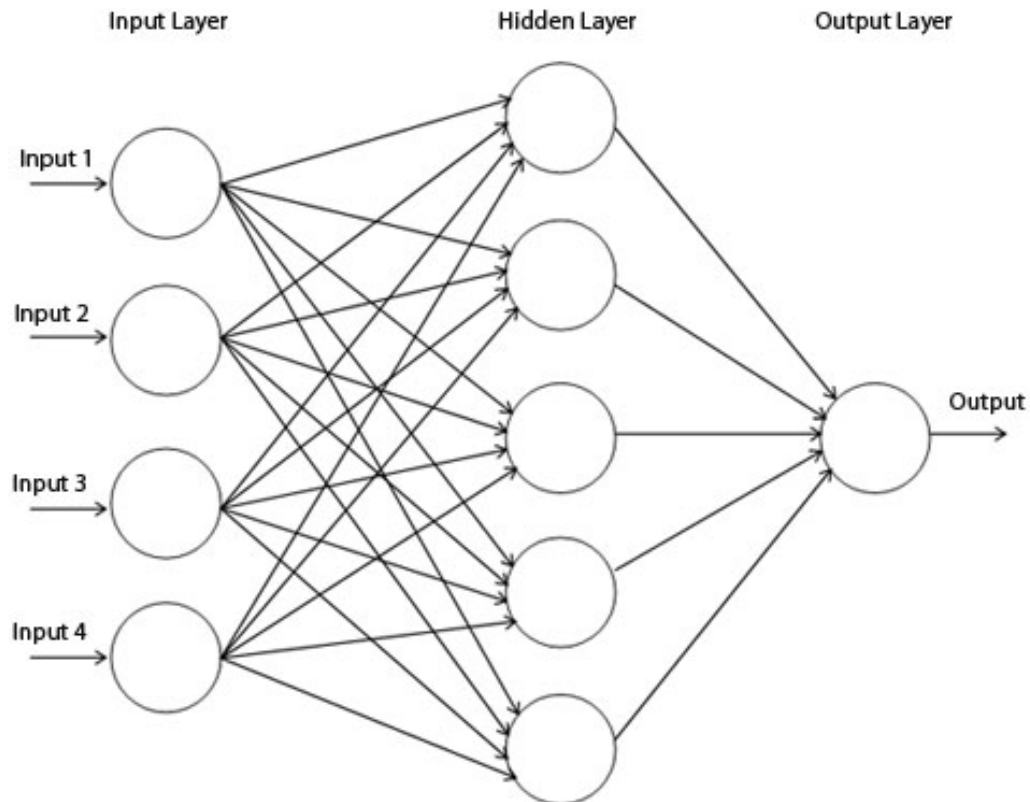


Figura 1.2: Struttura di un percettrone multistrato con un solo strato nascosto

Il *mlp* possiede le seguenti caratteristiche:

- ogni neurone è un percettore come quello descritto nella sezione 1.1. Ogni unità possiede quindi una propria funzione d'attivazione non lineare.
- a ogni connessione tra due neuroni corrisponde un peso sinaptico w
- è formato da 3 o più strati. In 1.2 è mostrato un MLP con uno strato di input; un solo strato nascosto (o *hidden layer*; ed uno di output.)
- l'uscita di ogni neurone dello strato precedente è l'ingresso per ogni neurone dello strato successivo. È quindi una rete *completamente connessa*. Tuttavia, si possono disconnettere selettivamente settando il peso sinaptico w a 0.
- la dimensione dell'input e la dimensione dell'output dipendono dal numero di neuroni di questi due strati. Il numero di neuroni dello strato nascosto è invece indipendente, anche se influenza di molto le capacità di apprendimento della rete.

Se ogni neurone utilizzasse una funzione lineare allora si potrebbe ridurre l'intera rete ad una composizione di funzioni lineari. Per questo - come detto prima - ogni neurone possiede una funzione di attivazione non lineare.

1.2.1 Strati Nascosti

I cosiddetti *hidden layers* sono una parte molto interessante della rete. Per il teorema di approssimazione universale [WApprox], **WApprox** [WConnessionismo] una rete con un singolo strato nascosto e un numero finito di neuroni, può essere addestrata per approssimare una qualsiasi funzione continua su uno spazio compatto di \mathbb{R}^n . In altre parole, un singolo strato nascosto è abbastanza potente da imparare un ampio numero di funzioni. Precisamente, una rete a 3 strati è in grado di separare regioni convesse con un numero di lati \leq numero neuroni nascosti.

Reti con un numero di strati nascosti maggiore di 3 vengono chiamate reti neurali profonde o *deep neural network*; esse sono in grado di separare regioni qualsiasi, quindi di approssimare praticamente qualsiasi funzione. Il primo e l'ultimo strato devono avere un numero di neuroni pari alla dimensione dello spazio di ingresso e quello di uscita. Queste sono le terminazioni della "*black box*" che rappresenta la funzione che vogliamo approssimare.

L'aggiunta di ulteriori strati non cambia *formalmente* il numero di funzioni che si possono approssimare; tuttavia vedremo che nella pratica un numero elevato di strati migliori di gran lunga le performance della rete su determinati task, essendo gli hidden layers gli strati dove la rete memorizza la propria rappresentazione astratta dei dati in ingresso. Nel capitolo 4 vedremo un'architettura all'avanguardia con addirittura 152 strati.

1.3 Caso di studio: prevedere il profitto di un ristorante

Prendendo spunto dalla traccia d'esame di Sistemi Intelligenti M del 2 Aprile 2009:

"Loris è figlio della titolare di una famoso spaccio di piadine nel Riminese e sta tornando in Italia dopo aver frequentato con successo un prestigioso Master in Business Administration ad Harvard, a cui si è iscritto inseguendo il sogno di esportare in tutto il mondo la piadina romagnola. Nel

lungo viaggio in prima classe, medita su come presentare alla mamma, che sa essere un tantino restia alle innovazioni, il progetto di aprire un ristorante a New York City."

Loris ha esportato con successo la piadina a NY, si veda [WGradisca] ma col passare degli anni ha notato alcuni problemi e vuole utilizzare di nuovo le sue brillanti capacità analitiche per migliorare il profitto del suo ambizioso ristorante.

I problemi sono 2:

1. il ristorante è conosciuto ormai - si sa che tutti vogliono mangiare italiano - ma il numero dei coperti era rimasto a 22, come quelli iniziali;
2. gli orari di apertura sono troppo lunghi e vi sono alcune zone morte dove il costo di mantenere aperto il ristorante è maggiore rispetto al ricavo dei pochi clienti che si siedono a mangiare durante quelle ore;

Secondo la National Restaurant Association [WProfit] [WRRG] il profitto medio lordo annuo di un ristorante negli Stati Uniti varia dal 2 al 6%. Così Loris ha collezionato alcuni dati riguardo gli ultimi anni e - attratto da tutta quest'entusiasmo attorno alle reti neurali - decide di provare ad utilizzarle per trovare il trade-off ottimale di coperti e di orari di apertura settimanali per massimizzare il profitto del suo ristorante.

1.4 Implementazione da zero di un MLP

Per il suddetto caso di studio si è implementato da zero un percettrone multistrato a 3 strati come quello illustrato in sezione 1.2. Vediamo qui di seguito le varie parti da implementare passo passo per costruire un MLP, addestrarlo e verificare che l'addestramento sia stato eseguito in maniera corretta. Il progetto è realizzato in Lua, per utilizzare il framework per il *Machine Learning* **Torch** e mantenere la consistenza con gli altri capitoli, nei quali si userà ancora Torch per addestrare reti neurali molto più complesse.

1.4.1 Dataset e Architettura

Nei vari anni, Loris ha cambiato le due variabili in gioco annotando di volta in volta i risultati. Siccome i coperti erano troppo pochi e le file d'attesa erano troppo lunghe, il ristorante perdeva alcuni clienti. Quindi Loris ha portato i coperti a 25 e diminuendo le ore settimanali a 38, ed i profitti sono aumentati. Tuttavia, non era raro che ancora qualche cliente dovesse aspettare in piedi per troppo tempo (si sa la vita a NY è frenetica), finendo poi per scegliere un ristorante adiacente. Inoltre, aveva diminuito le ore di troppo; nel weekend i clienti arrivavano fino a tardi, quindi rimanere aperti un'ora in più sarebbe stato lungimirante. Così, dopo l'allargamento della sala principale, ha aggiunto altri coperti ed aumentato le ore settimanali a 40, segnando un record personale di 4.4% di profitti annui.

Quindi, i dati in ingresso ed in uscita, in X e Y rispettivamente sono:

$$X = \begin{pmatrix} 22 & 42 \\ 25 & 38 \\ 30 & 40 \end{pmatrix} \quad Y = \begin{pmatrix} 2.8 \\ 3.4 \\ 4.4 \end{pmatrix}$$

Osservando le dimensioni dei dati si nota che la rete deve avere 2 input e dare in uscita 1 output, che chiameremo \hat{y} , in contrapposizione a y che è l'uscita desiderata. Per

quanto detto nella sezione 1.2, il MLP deve avere 2 neuroni nello strato di ingresso ed 1 solo in uscita. Inoltre, avrà uno strato nascosto con 3 neuroni. La dimensione di ogni strato fa parte di un insieme di parametri che viene deciso "a mano" sperimentando, i cosiddetti *hyperparameters*. Questi parametri non vengono aggiornati durante l'addestramento - come i pesi della rete - ma vengono decisi a priori. In figura 1.3 è mostrata l'architettura generale della nostra rete.

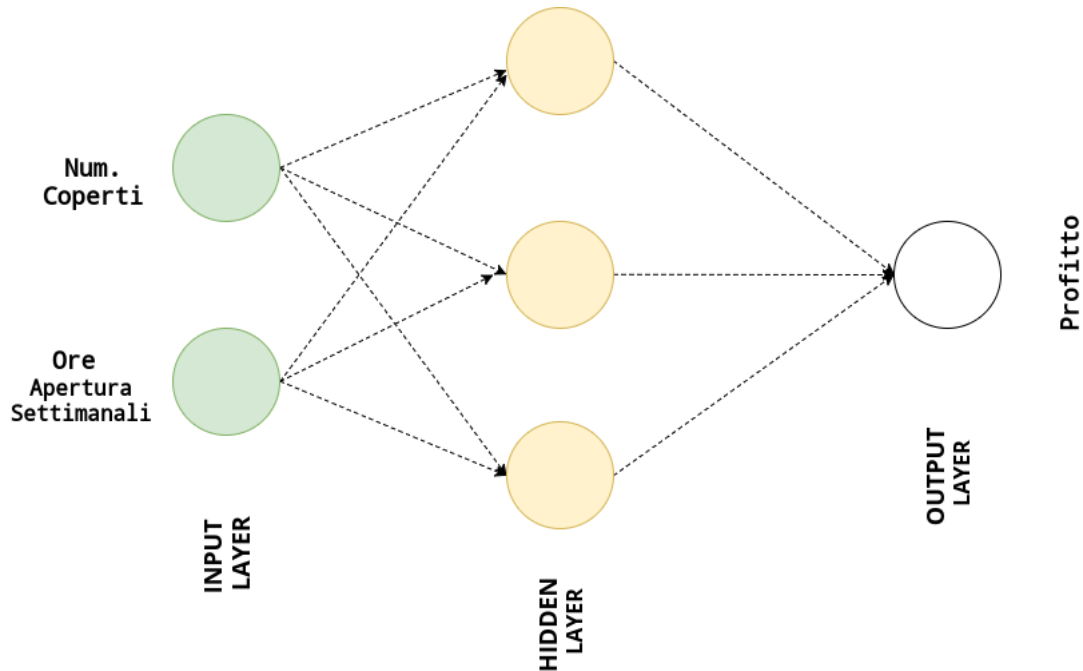


Figura 1.3: Architettura del perceptron multistrato per la previsione dei profitti

Di seguito gli snippet di codice per la definizione dei dati e dell'architettura della rete secondo lo schema appena presentato.

```

1 ----- Part 1 -----
2 th = require 'torch'
3 bestProfit = 6.0
4 -- X = (num coperti, ore di apertura settimanali), y = profitto lordo
   annuo in percentuale
5 torch.setdefaulttensortype('torch.DoubleTensor')
6 X = th.Tensor({{22,42}, {25,38}, {30,40}})
7 y = th.Tensor({{2.8}, {3.4}, {4.4}})
8
9 --normalize
10 normalizeTensorAlongCols(X)
11 y = y/bestProfit

```

```

1 ----- Part 2 -----
2 --creating the NN class in Lua, using a nice class utility
3 require 'class'
4
5 Neural_Network = class(function(net, inputs, hiddens, outputs)
6     net.inputLayerSize = inputs
7     net.hiddenLayerSize = hiddens
8     net.outputLayerSize = outputs
9     net.W1 = th.randn(net.inputLayerSize, net.hiddenLayerSize)
10    net.W2 = th.randn(net.hiddenLayerSize, net.outputLayerSize)
11 end)

```

1.4.2 Forward Propagation

Tabella 1.1: Variabili usate nel testo e nel codice

Variabili			
S. Codice	S. Matematico	Definizione	Dimensione
X	X	Esempi, 1 per riga	(numEsempi, inputLayerSize)
y	y	uscita desiderata	(numEsempi, outputLayerSize)
W1	$W^{(1)}$	Pesi layer 1	(inputLayerSize, hiddenLayerSize)
W2	$W^{(2)}$	Pesi layer 2	(hiddenLayerSize, outputLayerSize)
z2	$z^{(2)}$	Input layer 2	(numEsempi, hiddenLayerSize)
a2	$a^{(2)}$	Uscita layer 2	(numEsempi, hiddenLayerSize)
z3	$z^{(3)}$	Input layer 3	(numEsempi, outputLayerSize)

Nella tabella 1.1 sono elencate le variabili della rete. Gli input dei layer indicati con z possono anche essere chiamati "attività dei layer" (indicando l'attività sulle loro sinapsi); e $a^{(2)}$ indica l'uscita del neurone dopo aver applicato la sommatoria e la funzione di attivazione sulle attività provenienti dal layer precedente.

Per muovere i dati in parallelo attraverso la rete si usa la moltiplicazione fra matrici, per questo è molto comodo usare framework che supportano operazioni fra matrici come *Torch*, *Numpy* o *Matlab*. Per prima cosa, gli input del tensore X devono essere moltiplicati e sommati con i pesi del primo layer $W^{(1)}$, ottenendo l'ingresso per l'hidden layer:

$$z^{(2)} = XW^{(1)} \quad (1)$$

Si noti che $z^{(2)}$ è di dimensione 3x3, essendo X e $W^{(1)}$ di dimensione 3x2 e 2x3 rispettivamente.

Ora bisogna applicare la funzione di attivazione a $z^{(2)}$. Vi sono diverse funzioni di attivazione utilizzate per le reti neurali. Una delle prime a diventare popolare fu la funzione *sigmoide* [WSigmoid], utilizzata per questa rete. Vedremo nei capitoli successivi funzioni più efficaci.

$$a^{(2)} = f(z^{(2)}), \quad \text{ove } f = \text{sigmoide} \quad (2)$$

Per completare la *forward propagation*, bisogna seguire lo stesso procedimento per lo strato di output: sommare i contributi provenienti dall'hidden layer ed applicare la sigmoide:

$$z^{(3)} = a^{(2)}W^{(2)} \quad (3)$$

$$\hat{y} = f(z^{(3)}) \quad (4)$$

Essendo $a^{(2)}$ di dimensione 3×3 e $W^{(2)}$ 3×1 l'output \hat{y} sarà anch'esso di dimensione 3×1 , risultando quindi in una previsione per ogni esempio in ingresso.

Si noti come la moltiplicazione fra matrici renda tutto esprimibile in poche righe di codice.

```

1 --Note: I didn't implement manually the sigmoid function as Torch has one
  built-in.
2 --define a forward method
3 function Neural_Network:forward(X)
4     --Propagate inputs through network
5     self.z2 = th.mm(X, self.W1) --matrix multiplication
6     self.a2 = th.sigmoid(self.z2)
7     self.z3 = th.mm(self.a2, self.W2)
8     yHat = th.sigmoid(self.z3)
9     return yHat
10 end

```

1.4.3 Backpropagation

Come si fa ad addestrare una rete multistrato con diversi neuroni per strato, ognuno dei quali con uscita non lineare? Tramite l'algoritmo di *backpropagation of errors*, scoperto da Rumelhart-Hinton-Williams nel 1985.

Non si può introdurre l'algoritmo di "*backprop*" senza prima spiegare il concetto di funzione di costo (o *loss function*). Nelle reti neurali (e più specificatamente nell'apprendimento supervisionato [**WSupervised**]), la funzione di costo misura la discrepanza tra l'uscita desiderata e l'effettivo output della rete. È quindi una misura dell'errore della rete, per cui l'obiettivo dell'apprendimento è diminuire questa funzione. Come per la funzione di attivazione, anche in questo caso ci sono ampie possibilità di scelta [**WLoss**] a seconda del task su cui la rete viene addestrata.

E di nuovo, come per la funzione di attivazione, si è scelta una delle funzione di costo più popolari: l'errore quadratico medio.

$$J = \sum_{j=1}^n \frac{1}{2} (y - \hat{y})^2 \quad (5)$$

La *backprop* ha alcuni requisiti:

- Reti stratificate
- Ingressi a valori reali $\in [0, 1]$
- Neuroni non lineari con funzione di uscita sigmoidale (o altra fz. di attivazione derivabile)

Sotto queste condizioni l'algoritmo sfrutta la regola della catena [**WChain**] per la derivazione di funzione composte, per calcolare il gradiente della *funzione di costo*. I pesi della rete vengono quindi aggiornati secondo la *discesa del gradiente* (figura 1.4); ovvero variano in maniera tale da minimizzare la funzione di costo J .

Viene chiamato *backward propagation of errors* poiché l'errore calcolato a partire dall'output della rete viene distribuito in maniera proporzionale all'indietro, su tutti i neuroni della rete. È importante quindi, spezzare il calcolo del gradiente dell'errore in derivate parziali, dall'ultimo strato fino al primo, e poi combinarle insieme.

Si noti che le equazioni (1-5) formano una un'unica equazione che lega J a $X, y, W^{(1)}, W^{(2)}$. Tenendo questo in mente, si applica la regola della catena.

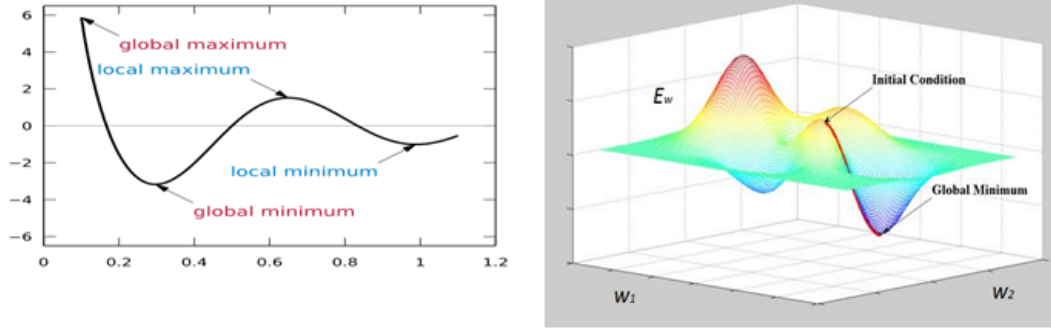


Figura 1.4: Cercare il minimo di una fz. seguendo la discesa del gradiente

Partendo dal fattore riguardante lo strato di output si ha:

$$\frac{\partial J}{\partial W^{(2)}} = \sum \frac{\partial \frac{1}{2}(y - \hat{y})^2}{\partial W^{(2)}}$$

Sviluppando i calcoli si ottiene:

$$\frac{\partial J}{\partial W^{(2)}} = -(y - \hat{y}) \frac{\partial \hat{y}}{\partial W^{(2)}}$$

L'equazione (4) indica che \hat{y} è la funzione di attivazione di $z^{(3)}$. Da cui:

$$\frac{\partial J}{\partial W^{(2)}} = -(y - \hat{y}) \frac{\partial \hat{y}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial W^{(2)}}$$

Il 2 membro dell'equazione è semplicemente la derivata della funzione di attivazione sigmoide:

$$f(z) = \frac{1}{1 + e^{-z}}$$

$$f'(z) = \frac{e^{-z}}{(1 + e^{-z})^2}$$

Definiamo quindi la derivata della funzione sigmoid in Lua:

Quindi l'equazione così ottenuta è:

$$\frac{\partial J}{\partial W^{(2)}} = -(y - \hat{y}) f'(z^{(3)}) \frac{\partial z^{(3)}}{\partial W^{(2)}}$$

Infine, dobbiamo trovare $\frac{\partial z^{(3)}}{\partial W^{(2)}}$, che rappresenta la variazione dell'attività del terzo layer rispetto ai pesi del secondo layer. Richiamando l'equazione (3):

$$z^{(3)} = a^{(2)} W^{(2)}$$

1.4.4 Verifica numerica del gradiente

1.4.5 Addestramento

1.5 Ottimizzazione: diverse tecniche

1.6 Overfitting: come rilevarlo e risolverlo

1.7 Risultati

Capitolo 2

Reti Neurali Convoluzionali

2.1 Breve introduzione

Le reti neurali convoluzionali, alle quali ci riferiremo con l'abbreviazione *CNN* - dall'inglese *Convolutional Neural Network*, sono un'evoluzione delle normali reti artificiali caratterizzate da una particolare architettura che le ha rese negli anni molto efficaci su diversi compiti.

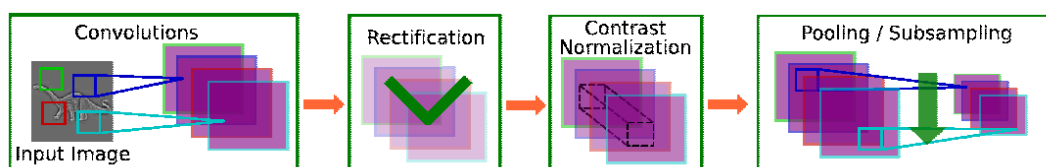


Figura 2.1: I diversi strati tipici di una CNN

2.2 Architettura

2.2.1 Strato di Convoluzione

2.2.2 Strato di Pooling

2.2.3 Strato completamente connesso (FC)

2.3 Applicazioni e risultati

Capitolo 3

CNN: implementazione e addestramento

3.1 Il framework Torch

3.2 Modello di CNN basato su LeNet

3.2.1 Strato di Convoluzione

3.2.2 Strato di Pooling

3.2.3 Strato completamente connesso (FC)

3.3 Dataset d'esempio: CIFAR

3.4 Addestramento su CIFAR

Vedremo nel capitolo seguente un confronto dei risultati di una CNN con un'architettura allo stato dell'arte addestrata sullo stesso dataset.

Capitolo 4

Addestrare un modello allo stato dell'arte

4.1 ResNet

4.2 Implementazione di Resnet di Facebook

4.3 Addestramento su CIFAR

4.4 LeNet vs Resnet: confronto

Capitolo 5

Caso d'uso attuale: fine-tuning su dataset arbitrario

Un caso di studio attuale è quello di prendere una rete allo stato dell'arte, già addestrata per mesi su un dataset enorme come quello di *Imagenet* per utilizzarla su un dataset arbitrario a nostra scelta, per scopi industriali o personali.

5.1 Dataset

5.2 Fine-Tuning

5.3 ModelZoo

5.3.1 Fine-tuning su Resnet

Facebook ha reso pubbliche le sue implementazioni di Resnet su github a tutti, cosicché possano essere utilizzate per qualsivoglia esperimento. Ci sono vari modelli, più o meno potenti a seconda degli strati della rete (da 18 a 140).

Capitolo 6

Conclusioni

Appendice A

Appendice A: Codice aggiuntionale ed il framework Torch

In Lua manca il costrutto delle classi. Si può tuttavia crearle utilizzando tables e meta-tables. Per il progetto si è utilizzato una piccola libreria, di seguito riportata.

```

1  -- class.lua
2  -- Compatible with Lua 5.1 (not 5.0).
3  function class(base, init)
4      local c = {} -- a new class instance
5      if not init and type(base) == 'function' then
6          init = base
7          base = nil
8      elseif type(base) == 'table' then
9          -- our new class is a shallow copy of the base class!
10         for i,v in pairs(base) do
11             c[i] = v
12         end
13         c._base = base
14     end
15     -- the class will be the metatable for all its objects,
16     -- and they will look up their methods in it.
17     c.__index = c
18
19     -- expose a constructor which can be called by <classname>(<args>)
20     local mt = {}
21     mt.__call = function(class_tbl, ...)
22         local obj = {}
23         setmetatable(obj, c)
24         if init then
25             init(obj, ...)
26         else
27             -- make sure that any stuff from the base class is initialized!
28             if base and base.init then
29                 base.init(obj, ...)
30             end
31         end
32         return obj
33     end
34     c.init = init
35     c.is_a = function(self, klass)
36         local m = getmetatable(self)
37         while m do
38             if m == klass then return true end
39             m = m._base
40         end
41         return false
42     end
43     setmetatable(c, mt)
44     return c

```

45 **end**